

Catherine

Regression and Estimation

slides 03 - regression
slides 04 - MLE

Linear Models (denominator layout)

Input: vector $x \in \mathbb{R}^D$ \longrightarrow Output: $y \in \mathbb{R}$

Data: $\langle (x_i, y_i) \rangle_{i=1}^N$

$$\text{Linear Model: } y = \underbrace{w_0 + w_1 x_1 + \dots + w_D x_D}_{\text{bias / intercept}} + \underbrace{\varepsilon}_{\text{noise / uncertainty}} = w^T x + \varepsilon$$

$$\hat{y} \in \mathbb{R}^{N \times 1} = X \in \mathbb{R}^{N \times (D+1)} \quad w \in \mathbb{R}^{(D+1) \times 1}$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_N^T \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix} = \begin{bmatrix} x_{10} & \cdots & x_{1D} \\ x_{20} & \cdots & x_{2D} \\ \vdots & \ddots & \vdots \\ x_{N0} & \cdots & x_{ND} \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix}$$

training phase: estimate \mathbf{w} from data

testing phrase: predict $\hat{y}_{\text{new}} = \hat{w}^T x_{\text{new}}$

Method: minimize average squared error $\frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$

$$\begin{aligned}
 \text{Loss function: } \mathcal{L}(w) &= \frac{1}{2N} \sum_{i=1}^N (x_i^T w - y_i)^2 = \frac{1}{2N} (Xw - y)^T (Xw - y) \\
 &= \frac{1}{2N} [w^T (X^T X) w - w^T X^T y - y^T X w + y^T y] \\
 &= \frac{1}{2N} [w^T (X^T X) w - 2y^T X w + y^T y]
 \end{aligned}$$

$$\nabla_w \mathcal{L} = \frac{1}{N} [(X^T X) w - X^T y]$$

By setting $\nabla_w \mathcal{L} = 0$, we get: $w = \underbrace{(X^T X)^{-1}}_{\text{Hessian}} X^T y$

Assumptions for unbiasedness (OLS estimators)

1. Linear in parameters: linear in ω
 2. Random sampling: (X_i, Y_i) iid, $E(Y_i|X) = E(Y_i|X_i)$, (not necessary)
 3. No perfect collinearity: $X^T X$ invertible, or $\text{rank}(X^T X) \xrightarrow{P=1} D+1$
 4. Zero conditional mean: $E(\varepsilon_i | X_i) = 0$
mean independence: $E(\varepsilon_i) = 0, E(\varepsilon_i | X_j) = 0, E(\varepsilon_i | X_i) = 0$ (if independent)
strictly exogenous: $E(\varepsilon_i | X)$, contemporaneously exogenous: $E(\varepsilon_i | X_i)$
 5. Gauss-Markov Assumptions). Conditional Homoskedasticity: $\text{Var}(\varepsilon_i | X)^2 = \sigma^2$
 $E(\varepsilon_i \varepsilon_j | X) = 0$ (conditional spatial / serial uncorrelations) $\Rightarrow \text{Cov}(\varepsilon_i, \varepsilon_j) = 0$

More on OLS

Linear Model: $\hat{y} = w_0 + w_1 x_1 + \dots + w_D x_D + \varepsilon = w^T x + \varepsilon$

$$\hat{y} \in \mathbb{R}^{N \times 1} = X \in \mathbb{R}^{N \times (D+1)} \quad w \in \mathbb{R}^{(D+1) \times 1}, \quad k=D+1$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_N^T \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix} = \begin{bmatrix} x_{10} & \dots & x_{1D} \\ x_{20} & \dots & x_{2D} \\ \vdots & \ddots & \vdots \\ x_{N0} & \dots & x_{ND} \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_D \end{bmatrix}$$

$$\mathbb{R}^{D+1} \xrightarrow[N]{\text{sample matrix}} \mathbb{R}^N$$

$$R(X^T)$$

$$\begin{aligned} X \hat{w} &= X(X^T X)^{-1} X^T Y \\ &\quad \text{projection matrix: } P = P^2, P X = X \\ R(X) &\rightarrow Y \\ N(X^T) &\rightarrow \hat{\varepsilon} = (I - P) Y = M Y \\ &= M(Xw + \varepsilon) = M\varepsilon \\ &= (I - X(X^T X)^{-1} X^T) \varepsilon \\ &\quad \text{symmetric and idempotent} \\ SSR &= \hat{\varepsilon}' \hat{\varepsilon} = \varepsilon' M\varepsilon \\ &= Y M Y \end{aligned}$$

Estimation of $\sigma^2 (= E(\varepsilon_i^2))$:

1. Method of Moments (MoM): $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{1}{n} \hat{\varepsilon}' \hat{\varepsilon}$

2. Unbiased Estimator: $S^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{SSR}{n-k}$

Variance Decomposition:

$$Y^T Y = \hat{Y}^T \hat{Y} + \hat{\varepsilon}^T \hat{\varepsilon} \Rightarrow SST = SSE + SSR$$

$$R^2 = 1 - \frac{SSR}{SST} = 1 - \frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad \bar{R}^2 = 1 - \frac{SSR / (n-k)}{SST / (n-1)} = 1 - \frac{n-1}{n-k} (1 - R^2)$$

Proof of unbiasedness: $\hat{w} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (Xw + \varepsilon) = w + (X^T X)^{-1} X^T \varepsilon$

$$E(\hat{w}|X) = w + E((X^T X)^{-1} X^T \varepsilon | X) = w + (X^T X)^{-1} X^T E(\varepsilon | X) = w$$

$$\begin{aligned} E(S^2 | X) &= \frac{1}{n-k} E(\varepsilon^T M\varepsilon | X) = \frac{1}{n-k} E[\text{tr}(\varepsilon^T M\varepsilon) | X] \\ &= \frac{1}{n-k} E[\text{tr}(M\varepsilon'\varepsilon | X)] = \frac{1}{n-k} \text{tr}[M E(\varepsilon'\varepsilon | X)] = \frac{1}{n-k} \text{tr}(M) = \sigma^2 \end{aligned}$$

Variance of \hat{w} : $\text{Var}(\hat{w}|X) = \text{Var}((X^T X)^{-1} X^T \varepsilon | X) = (X^T X)^{-1} X^T \text{Var}(\varepsilon | X) [(X^T X)^{-1} X^T]^T$

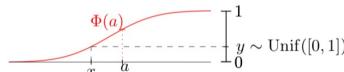
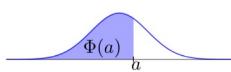
$$= (X^T X)^{-1} X^T \sigma^2 I_n X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

$$= \frac{Y^2}{SST_x} (D=1), \quad \frac{SST_j (1-R_j^2)}{SST_x (1-R_j^2)} (D>1)$$

$$GLR: \hat{w} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \left(\frac{\text{Cov}(X, Y)}{\text{Var}(X)} \right)$$

Probability Review

univariate normal distribution: $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$, $X \sim N(\mu, \sigma^2)$
 inverse transform sampling: $y = \Phi^{-1}(z; 0, 1)$, $Z \sim \text{Unif}([0, 1])$



bivariate normal distribution: $X_1 \sim N(\mu_1, \sigma_1^2)$, $X_2 \sim N(\mu_2, \sigma_2^2)$

$$\begin{aligned} f(x_1, x_2) &= \frac{1}{2\pi|\Sigma|^{1/2}} \exp\left(-\left(\frac{(x_1-\mu_1)^2}{2\sigma_1^2} + \frac{(x_2-\mu_2)^2}{2\sigma_2^2}\right)\right) \\ &= \frac{1}{2\pi|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right) \\ \text{where } \Sigma &= \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}, \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \end{aligned}$$

Covariance Matrix:

$$\text{Cov}(X) = E[(X - E(X))(X - E(X))^T] = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_D) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_D) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_D, X_1) & \text{Cov}(X_D, X_2) & \dots & \text{Var}(X_D) \end{bmatrix}$$

$$N(\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right), \Sigma = \text{Cov}(X)$$

Eigendecomposition: $\Sigma = Q \Lambda Q^{-1}$
 eigenvectors $\overset{\curvearrowleft}{\underset{\curvearrowright}{\sim}}$ eigenvalues

the eigenvector with the largest eigenvalue corresponds to the direction with the greatest variance and vice-versa

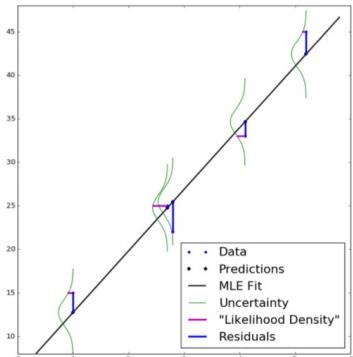
Maximum Likelihood

Linear Model: $y = w_0 x_0 + w_1 x_1 + \dots + w_D x_D + \varepsilon = w^T x + \varepsilon$
 aleatoric uncertainty

prediction: $E[y | x, w] = f_w(x) = w^T x$

$p(y | x, w) = N(y | w^T x, \sigma^2)$, $y \sim w^T x + N(0, \sigma^2)$
 (alternatively, $\varepsilon \sim N(0, \sigma^2)$, Gaussian Noise)

Likelihood of Linear Regression (Gaussian Noise Model)



MLE Estimator

Find parameters which maximise the likelihood.

(product of "likelihood density" — segments)

Least Square Estimator

Find parameters which minimise the sum of squares of the residuals

(sum of squares of the $|I|$ segments).

$$\begin{aligned}
 p(y_1, y_2, \dots, y_N | X_1, X_2, \dots, X_N, w, \tau) &= \prod_{i=1}^N p(y_i | X_i, w, \tau), \text{ model parameters} \\
 &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{(y_i - w^T X_i)^2}{2\tau^2}\right) \\
 &= \left(\frac{1}{2\pi\tau^2}\right)^{N/2} \exp\left(-\frac{1}{2\tau^2} \sum_{i=1}^N (y_i - w^T X_i)^2\right)
 \end{aligned}$$

$$\text{log-likelihood: } LL(y_1, y_2, \dots, y_N | X_1, X_2, \dots, X_N, w, \tau)$$

$$= -\frac{N}{2} \log(2\pi\tau^2) - \frac{1}{2\tau^2} \sum_{i=1}^N (y_i - w^T X_i)^2$$

$$\Rightarrow LL(y | X, w, \tau) = -\frac{N}{2} \log(2\pi\tau^2) - \frac{1}{2\tau^2} (Xw - y)^T (Xw - y)$$

$$NLL(y | X, w, \tau) = \frac{N}{2} \log(2\pi\tau^2) + \frac{1}{2\tau^2} (Xw - y)^T (Xw - y)$$

minimize NLL $\Rightarrow w, \tau^2 \Leftrightarrow$ least square method

$$\Rightarrow \begin{cases} w_{ML} = (X^T X)^{-1} X^T y \\ \tau_{ML}^2 = \frac{1}{N} (Xw_{ML} - y)^T (Xw_{ML} - y) \end{cases}$$

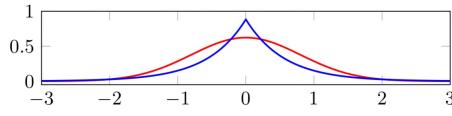
Outliers and Laplace Distribution

$$\text{Linear Model: } y = w_0 x_0 + w_1 x_1 + \dots + w_D x_D + \varepsilon = w \cdot x + \varepsilon$$

$$\varepsilon \sim \text{Lap}(0, b), \text{ pdf}(\varepsilon) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

$$\begin{aligned}
 p(y_1, \dots, y_N | X_1, \dots, X_N, w, b) \\
 = \frac{1}{(2b)^N} \exp\left(-\frac{1}{b} \sum_{i=1}^N |y_i - w^T X_i|\right)
 \end{aligned}$$

$$NLL(y | X, w, b) = \frac{1}{b} \sum_{i=1}^N |y_i - w^T X_i| + N \log(2b)$$



Fitting, Regularization and Bayesian

Slides 05 - fitting
slides 06 - regularization
slides 07 - Bayesian

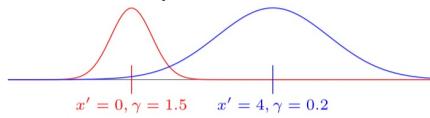
Basis Function Expansion

polynomial basis expansion: $y = w \cdot \phi(X) + \epsilon$

{ linear model: $\phi(X) = [1, x_1, x_2]$

quadratic model: $\phi(X) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$

radial basis function (RBF) kernel: $K(x', x) = \exp(-\gamma \|x - x'\|^2)$



$K(x', x) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} = \phi(x') \cdot \phi(x)$

polynomial kernel: $K(x', x) = (x^T x + c)^\alpha$

{ underfitting: γ too small (large width)

overfitting: γ too large (small width)

1. choose centres: $\mu_1, \mu_2, \dots, \mu_m$

2. feature map: $\phi(X) = [1, K(\mu_1, X), \dots, K(\mu_m, X)]$

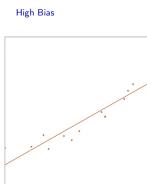
3. model: $y = w_0 + w_1 K(\mu_1, X) + \dots + w_m K(\mu_m, X) + \epsilon = w \cdot \phi(X) + \epsilon$

curse of dimensionality: might need exponentially large (in the dimension) sample for using modest width kernels.

The Bias Variance Tradeoff

{ underfitting: high bias (high similar errors)

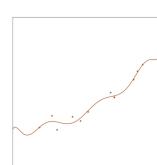
{ overfitting: high variance (training \approx test errors)



Learning Curves:

training set + validation + test

plot errors as a function of training data size



Regularization

Ridge Regression Objective $L_{\text{ridge}}(w) = (Xw - y)^T (Xw - y) + \lambda w^T w$
(*) before optimization, standardize all inputs (mean 0 and variance 1)

$$\begin{aligned} L_{\text{ridge}}(w) &= (Xw - y)^T (Xw - y) + \lambda w^T w \\ &= w^T (X^T X) w - 2y^T Xw + y^T y + \lambda w^T w \end{aligned}$$

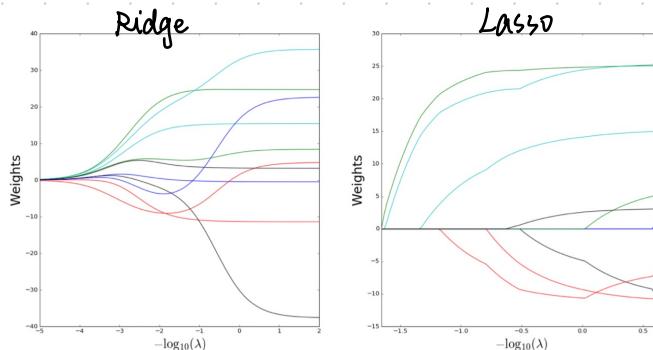
L_2 -regularization
(weight-decay)

$$\begin{aligned}\nabla_w \mathcal{L}_{\text{ridge}} &= 2(X^T X) w - 2X^T y + 2\lambda w \\ &= 2[(X^T X + \lambda I_D) w - X^T y] \\ \Rightarrow w_{\text{ridge}} &= (X^T X + \lambda I_D)^{-1} X^T y\end{aligned}$$

Lasso Regression Objective $\mathcal{L}_{\text{ridge}}(w) = (Xw - y)^T (Xw - y) + \lambda \sum_{i=1}^D |w_i|$
(*) before optimization, standardize all inputs (mean 0 and variance 1)

$$\begin{aligned}\mathcal{L}_{\text{ridge}}(w) &= (Xw - y)^T (Xw - y) + \lambda w^T w \\ &= w^T (X^T X) w - 2y^T X w + y^T y + \lambda w^T w\end{aligned}$$

l_1 -regularization
(Least Absolute Shrinkage and Selection Operator)



when using Lasso:
weights are often exactly 0
 \Rightarrow sparse models

Feature Selection

Forward Search:

$\Theta(n^2)$

- Set set of selected features to $F := \emptyset$
- Repeat the following until $F = \{1, 2, \dots, n\}$:
 - set $F_i := F \cup \{i\}$ for $i \in \{1, 2, \dots, n\} \setminus F$
 - evaluate generalization error when using only features from F_i
 - set new F to the best feature subset found
- Return best overall feature subset found

Filter feature selection:

criterion (mutual information).

$$I(X, Y) = \sum_x \sum_y p(X=x, Y=y) \cdot \log \frac{p(X=x, Y=y)}{p(X=x) \cdot p(Y=y)}$$

compute mutual information for all features
retain top k features

Exercises

1.

$$\mathbb{E}_{\mathcal{D}} \left[\|\hat{\mathbf{w}}(\mathcal{D}) - \mathbf{w}^*\|^2 \right] = \mathbb{E}_{\mathcal{D}} \left[\|\hat{\mathbf{w}}(\mathcal{D})\| \right]^2 + \mathbb{E}_{\mathcal{D}} \left[\|\hat{\mathbf{w}}(\mathcal{D}) - \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})]\|^2 \right]$$

MSE Bias² Variance

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\|\hat{\mathbf{w}}(\mathcal{D}) - \mathbf{w}^*\|^2] &= \mathbb{E}_{\mathcal{D}} [\|\hat{\mathbf{w}}(\mathcal{D}) - \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})] + \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})] - \mathbf{w}^*\|^2] \\ &= \mathbb{E}_{\mathcal{D}} [\|\hat{\mathbf{w}}(\mathcal{D}) - \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})]\|^2] + \mathbb{E}_{\mathcal{D}} [\|\mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})] - \mathbf{w}^*\|^2] \\ &\quad + 2 \mathbb{E}_{\mathcal{D}} [(\hat{\mathbf{w}}(\mathcal{D}) - \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})]) \cdot (\mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})] - \mathbf{w}^*)] \\ &= \mathbb{E}_{\mathcal{D}} [\|\hat{\mathbf{w}}(\mathcal{D}) - \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})]\|^2] + \|\mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})] - \mathbf{w}^*\|^2 \\ &\quad \leftarrow \mathbb{E}_{\mathcal{D}} (\hat{\mathbf{w}}(\mathcal{D}) - \mathbb{E}_{\mathcal{D}} [\hat{\mathbf{w}}(\mathcal{D})]) = \mathbb{E}_{\mathcal{D}} (\hat{\mathbf{w}}(\mathcal{D})) - \mathbb{E}_{\mathcal{D}} [\mathbb{E}_{\mathcal{D}} (\hat{\mathbf{w}}(\mathcal{D}))] = 0 \\ &\qquad\qquad\qquad \text{unrelated to } \mathcal{D}, \\ &\qquad\qquad\qquad \text{thus constant wrt } \mathbb{E}_{\mathcal{D}} \cdot \end{aligned}$$

2. Huber Loss

In this question, we will investigate the *Huber loss* in a linear regression setting. Given arbitrary but fixed parameters $\lambda, \mu \in \mathbb{R}$ such that $\lambda, \mu > 0$, the Huber loss is given by the function $h_{\lambda, \mu} : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$h_{\lambda, \mu}(z) = \begin{cases} \lambda \left(|z| - \frac{\lambda}{4\mu} \right) & \text{if } |z| \geq \frac{\lambda}{2\mu}, \\ \mu z^2 & \text{otherwise.} \end{cases} \quad (\text{regime}).$$

Given a vector $\mathbf{z} = (z_1, \dots, z_D) \in \mathbb{R}^D$, we extend $h_{\lambda, \mu}$ such that $h_{\lambda, \mu}(\mathbf{z}) = \sum_{i=1}^D h_{\lambda, \mu}(z_i)$. Recall that when dealing with absolute values, the sign function defined as follows is often helpful:

$$\text{sign}(z) = \begin{cases} 1 & \text{if } z > 0, \\ -1 & \text{otherwise.} \end{cases}$$

$$\mathcal{H}(\mathbf{z}; \mathcal{D}) = h_{\lambda, \mu}(\mathbf{z}) + \frac{1}{N} \sum_{i=1}^N \ell(y_i - \mathbf{z}^\top \cdot \mathbf{x}_i),$$

$$\mathcal{S}(\mathbf{v}, \mathbf{w}; \mathcal{D}) = \lambda \|\mathbf{v}\|_1 + \mu \|\mathbf{w}\|_2^2 + \frac{1}{N} \sum_{i=1}^N \ell(y_i - (\mathbf{v} + \mathbf{w})^\top \cdot \mathbf{x}_i).$$

Sklearn Python Programming

```
from sklearn.pipeline import make_pipeline  
from sklearn.preprocessing import PolynomialFeatures, StandardScaler  
from sklearn.linear_model import Ridge, Lasso  
from sklearn.metrics import mean_squared_error  
  
degree = 2  
ridge_model = make_pipeline(PolynomialFeatures(degree),  
                           StandardScaler(), Ridge(alpha=lamb),  
                           verbose=True)  
ridge_model.fit(X_train, y_train)  
ridge_validation_pred = ridge_model.predict(X_validation)  
ridge_mse = mean_squared_error(y_validation, ridge_validation_pred)
```

} First polynomialize
Then standardize

```
# model selection  
from sklearn.model_selection import GridSearchCV
```

```
param_grid_ridge = {  
    'polynomialfeatures_degree': [2, 3, 4],  
    'ridge_alpha': [0.01, 0.1, 1, 10, 100]  
}  
ridge_model = make_pipeline(PolynomialFeatures(),  
                           StandardScaler(), Ridge())  
ridge_cv = GridSearchCV(ridge_model, param_grid_ridge, cv=5,  
                       scoring="neg_mean_squared_error")  
ridge_model.fit(X_train, y_train)  
best_ridge_model = ridge_cv.best_estimator_  
ridge_test_pred = ridge_model.predict(X_test)  
ridge_mse = mean_squared_error(y_test, ridge_test_pred)
```

dictionary key
corresponds to
preprocessing
classes

K-fold Cross Validation

	(K=5)				
Run 1	train	train	train	train	valid
Run 2	train	train	train	valid	train
Run 3	train	train	valid	train	train
Run 4	train	valid	train	train	train
Run 5	valid	train	train	train	train

when K=N, LOOCV (Leave One Out Cross Validation)

Other Tips for programming

```
# bar chart  
v = []  
for i in range(3, 10):  
    v.append(y_train[y_train == i].size)  
plot.bar(range(3, 10), v)
```

```
# numpy  
means = np.mean(X_train, axis=0) (vertically)  
stds = np.std(X_train, axis=0)  
X_train_n = (X_train - means) / stds (horizontal)  
X_train_n = np.hstack((X_train_n, np.full((N_train, 1), 1))) shape ; fill value  
W = np.linalg.inv(X_train_n.transpose().dot(X_train_n)) tuple  
.dot(X_train_n.transpose()).dot(y_train)  
y_hat_train = X_train_n.dot(W)
```

$$X_{\text{train}} = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \quad \vdots \quad \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \quad N \times D$$

Bayesian Approach

Frequentists: fixed parameters, $P(D|w)$

Bayesians: fixed data, $p(w|D)$

Aleatoric uncertainty: stochastic nature of variables (more data useless)

Epistemic uncertainty: about choice of model (reducible with more data)

$$\text{Bayes' Theorem: } \underbrace{p(w|D)}_{\text{posterior}} = \frac{p(D|w) \cdot p(w)}{p(D)} \quad \begin{matrix} w \rightarrow D \\ \vdots \end{matrix}$$

$$\Rightarrow p(w|D) = \frac{p(y|w, X) \cdot p(w)}{p(D)}$$

prior: $p(\mu_s) = N(\mu_{\text{obs}}, \tau_{\text{obs}}^2)$

posterior: $p(\mu_s|D) = N(\mu_{N,s}, \tau_{N,s}^2)$ from samples $t_i \sim N(\mu, \tau^2)$

$$\text{, where: } \mu_{N,s} = \frac{1}{\frac{\tau^2}{\tau_{\text{obs}}^2} + N} \left[\frac{\tau^2}{\tau_{\text{obs}}^2} \mu_{\text{obs}} + \sum_{i=1}^N t_i \right], \tau_{N,s}^2 = \left(\frac{1}{\frac{\tau^2}{\tau_{\text{obs}}^2}} + \frac{N}{\tau^2} \right)^{-1}$$

Bernstein-von Mises Theorem:

under certain regularity assumptions, as $N \rightarrow \infty$, posterior converges to a normal distribution $N(\hat{w}_{\text{MLE}}, N^{-1} I_F(\hat{w}_{\text{MLE}})^{-1})$

centred on maximum likelihood estimate (\hat{w}_{MLE}) where $I_F(\hat{w}_{\text{MLE}})$ is the Fisher information matrix at \hat{w}_{MLE}

$$\Rightarrow p(t|s, D) = \int_w \underbrace{p(t|s, w)}_{\text{model}} \cdot \underbrace{p(w|D)}_{\text{posterior}} dw = N(\mu_{N,s}, \tau^2 + \tau_{N,s}^2)$$

$$E[t|s, D] = \int t \cdot p(t|s, D) dt = \mu_{N,s}$$

$$\text{pdf}(y|X, w) = \frac{1}{(2\pi\tau^2)^{N/2}} \prod_{i=1}^N \exp\left(-\frac{(y_i - w^T X_i)^2}{2\tau^2}\right) \Rightarrow \exp\left(-\frac{1}{2\tau^2}(y - Xw)^T(y - Xw)\right)$$

Gaussian prior: $p(w) = N(w|0, \Lambda) \propto \exp(-\frac{1}{2} w^T \Lambda^{-1} w)$

posterior: $\text{pdf}(w|D) \propto \underbrace{\text{likelihoood}}_{\text{likelihood}} \cdot \underbrace{p(w|0, \Lambda)}_{\text{prior}} = N(w|\hat{w}_N, \Sigma_N)$

$$\text{, where } \Sigma_N = (\Lambda^{-1} + \frac{1}{\tau^2} X^T X)^{-1}, \hat{w}_N = \frac{1}{\tau^2} \sum_N X^T y$$

} complexity!

$$\text{pdf}(y|D, X_{\text{new}}, \tau^2) = \int_w \text{pdf}(y|X_{\text{new}}, w) \text{pdf}(w|D) dw$$

$$= N(y | \hat{w}_N^T X_{\text{new}}, \tau^2 + X_{\text{new}}^T \Sigma_N X_{\text{new}})$$

aleatoric epistemic

Maximum a Posteriori (MAP)

prior over w takes the more general form: $\text{pdf}(w) \propto \exp(-R(w))$, $R: \mathbb{R}^p \rightarrow \mathbb{R}$

$$w_{\text{MAP}} \in \arg_w \max \text{pdf}(w | D) = \arg_w \max \log \text{pdf}(w | D)$$
$$\Rightarrow \text{pdf}(w | D) \propto \exp\left(-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)\right) \cdot \exp(-R(w))$$
$$= \exp\left(-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) - R(w)\right)$$

the MAP loss function:

$$L_{\text{MAP}}(w) \propto -\log(\text{pdf}(w | D)) = \underbrace{\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)}_{\text{least squares}} - \underbrace{R(w)}_{\text{regularization}}$$

$$L_{\text{ridge}}(w) \propto \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) + \frac{\lambda}{2\sigma^2} w^T w = L_{\text{MAP}}(w)$$

the prior corresponds to $\exp(-\frac{\lambda}{2\sigma^2} w^T w) \Rightarrow N(0, \text{diag}(\frac{\sigma^2}{\lambda}))$

$$L_{\text{lasso}}(w) \propto \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) + \frac{\lambda}{2\sigma^2} \sum_{i=1}^D |w_i| = L_{\text{MAP}}(w)$$

the prior corresponds to $\exp(-\frac{\lambda}{2\sigma^2} \sum_{i=1}^D |w_i|) \Rightarrow \text{pdf}(w) = \prod_{i=1}^D \text{Lap}(0, \frac{2\sigma^2}{\lambda})$

$$L_{\text{me}}(w) \propto \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) + \lim_{\lambda \rightarrow 0} \frac{1}{2\sigma^2} w^T w = L_{\text{MAP}}(w)$$

the prior corresponds to $\exp(-\lim_{\lambda \rightarrow 0} \frac{1}{2\sigma^2} w^T w) \Rightarrow N(0, \underbrace{\lim_{\lambda \rightarrow 0} \text{diag}(\frac{\sigma^2}{\lambda})}_{\substack{\text{improper prior} \\ \Rightarrow \text{entire } \mathbb{R}^p})$

Optimization

slides 08 - optimization
slides 09 - optimization 2

Convex Optimization

Definition: the entire set $\mathbb{R}^D : \lambda x + (1-\lambda)y \in \mathbb{R}^D, \forall x, y \in \mathbb{R}^D$
function $f: \mathbb{R}^D \rightarrow \mathbb{R} : f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$

Theorem: the set of positive semi-definite matrices is convex
positive semi-definite. $A = A^T, x^T A x \geq 0$ for $\forall x \in \mathbb{R}^D$
set $S_+^D = \text{positive semidefinite cone}$
 $\forall A, B \in S_+^D, x^T (\lambda A + (1-\lambda)B) \cdot x = \lambda \cdot x^T A x + (1-\lambda) x^T B x \geq 0$

Examples for concave functions:

1. Affine functions: $f(x) = b^T x + c$
2. Quadratic functions: $f(x) = \frac{1}{2} x^T A x + b^T x + c$,
 A is symmetric positive semidefinite
3. Norms, in particular L^p -norms
4. Nonnegative weighted sums of convex functions

Theorem: For a convex optimization problem, all locally optimal points are globally optimal

Convex Optimization

Given convex functions f, g_1, \dots, g_m and affine functions h_1, \dots, h_n .

minimize $f(x)$
subject to $g_i(x) \leq 0, i \in \{1, 2, \dots, m\} \Rightarrow$ optimal value
 $h_j(x) = 0, j \in \{1, 2, \dots, n\} \Rightarrow$ optimal point

1. Linear Programming

minimize $c^T x + d$ s.t. $Ax \leq e, Bx = f$

2. Quadratically Constrained Quadratic Programming

minimize $\frac{1}{2} x^T B x + c^T x + d$ s.t. $\frac{1}{2} x^T Q_i x + r_i^T x + s_i \leq 0, Ax = b$

3. Semidefinite Programming

minimize $\text{tr}(C X)$ s.t. $\text{tr}(A_i X) = b_i, X$ positive semidefinite

Exercises

1. Optimization Method for ℓ_1 -regularization

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N |\mathbf{w}^\top \mathbf{x}_i - y_i| + \lambda \sum_{i=1}^D |w_i|$$

Consider a linear program with $2D+N$ variables:

$$w_1, \dots, w_D, k_1, \dots, k_D, \xi_1, \dots, \xi_N$$

The objective function can be expressed in:

$$\text{minimize } \sum_{i=1}^N \xi_i + \lambda \sum_{i=1}^D k_i$$

$$\text{subject to } w_i^\top \mathbf{x}_i - y_i \leq \xi_i, \quad i=1, 2, \dots, N$$

$$y_i - w_i^\top \mathbf{x}_i \leq \xi_i, \quad i=1, 2, \dots, N$$

$$w_i \leq k_i, \quad i=1, 2, \dots, D$$

$$-w_i \leq k_i, \quad i=1, 2, \dots, D$$

$$\mathcal{L}_{\text{lasso}}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \sum_{i=1}^D |w_i|$$

$$\text{Set } f(w) = \lambda \sum_{i=1}^D |w_i|$$

$$\nabla_w f(w) = \lambda \begin{bmatrix} \underset{\stackrel{+}{\text{if }}}{\text{sign}(w_1)} \cdot \mathbb{I}(w_1 \neq 0) \\ \vdots \\ \underset{\stackrel{+}{\text{if }}}{\text{sign}(w_D)} \cdot \mathbb{I}(w_D \neq 0) \end{bmatrix}, \text{ where } \mathbb{I}(w_i \neq 0) = \begin{cases} 1, & w_i \neq 0 \\ 0, & w_i = 0 \end{cases}$$

$$\text{Thus } \nabla_w \mathcal{L}_{\text{lasso}}(w) = 2(X^\top X w - X^\top y) + \lambda \begin{bmatrix} \underset{\stackrel{+}{\text{if }}}{\text{sign}(w_1)} \cdot \mathbb{I}(w_1 \neq 0) \\ \vdots \\ \underset{\stackrel{+}{\text{if }}}{\text{sign}(w_D)} \cdot \mathbb{I}(w_D \neq 0) \end{bmatrix}$$

sub-gradient descent update rule:

$$w_{t+1} = w_t - \eta (2(X^\top X w - X^\top y) + \lambda \begin{bmatrix} \underset{\stackrel{+}{\text{if }}}{\text{sign}(w_1)} \cdot \mathbb{I}(w_1 \neq 0) \\ \vdots \\ \underset{\stackrel{+}{\text{if }}}{\text{sign}(w_D)} \cdot \mathbb{I}(w_D \neq 0) \end{bmatrix})$$

Gradient Descent

Gradients: $\nabla_w f(w) = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \\ \frac{\partial f}{\partial w_n} \end{bmatrix}$

Hessians: $H = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_n} \\ \frac{\partial^2 f}{\partial w_2 \partial w_1} & \frac{\partial^2 f}{\partial w_2^2} & \cdots & \frac{\partial^2 f}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_n \partial w_1} & \frac{\partial^2 f}{\partial w_n \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_n^2} \end{bmatrix}$

Symmetric, capture the curvature of the surface

gradient descent: $w_{t+1} = w_t - \eta_t \nabla f(w_t)$, $\eta_t > 0$: learning rate / step size
gradient descent for least squares regression:

$$L(w) = (Xw - y)^T (Xw - y)$$

$$\nabla_w L = 2(X^T X w - X^T y)$$

complexity: $O(ND)$
($X^T X$ only calculate once)

Newton's Method

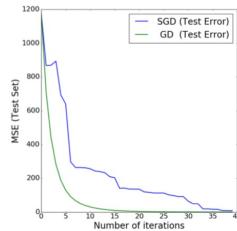
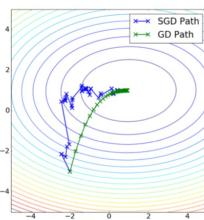
$$w_{t+1} = w_t - H_t^{-1} g_t \quad f_{\text{quad}}(w) = f(w_t) + g_t^T (w - w_t) + \frac{1}{2} (w - w_t)^T H_t (w - w_t)$$
$$\nabla_w f_{\text{quad}} = g_t + H_t (w - w_t) := 0$$

Sub-gradient Descent: range between the left and right derivatives

Stochastic Gradient Descent (SGD)

objective function: $L(w; D) = \frac{1}{N} \sum_{i=1}^N l(w; x_i, y_i) + \lambda R(w)$ regularization term
pick a random datapoint (x_i, y_i) and evaluate $g_i = \nabla_w l(w; x_i, y_i)$

$$E(g_i) = \frac{1}{N} \sum_{i=1}^N \nabla_w l(w; x_i, y_i)$$



Batch / Offline Learning:

$$w_{t+1} = w_t - \frac{1}{N} \sum_{i=1}^N \nabla_w l(w; x_i, y_i) - \lambda \nabla_w R(w)$$

Online Learning:

$$w_{t+1} = w_t - \eta \nabla_w l(w; x_t, y_t) - \lambda \nabla_w R(w)$$

Minibatch Online Learning

$$w_{t+1} = w_t - \frac{1}{b} \sum_{i=1}^b \nabla_w l(w; x_i, y_i) - \lambda \nabla_w R(w)$$

Other Optimization Techniques

First Order Methods / (Sub) Gradient Methods:

1. Nesterov's Accelerated Gradient
2. Line-Search to Find Step-Size
3. Momentum-based Methods
4. AdaGrad, AdaDelta, Adam, RMSprop

$$W_{t+1,i} \leftarrow W_{t,i} - \frac{\eta}{\sqrt{\sum_{s=1}^t g_{s,i}^2}} g_{t,i}$$

rare features can be most predictive
suitable in NLP

Second Order / Newton / Quasinewton Methods:

1. Conjugate Gradient Method
2. BFGS and L-BFGS

Classification

Slides 10 - classification (NBC)
Slides 11 - logistic regression
Slides 12 - SVM
Slides 13 - SVM (2)

Discriminative Model: distinguish decision boundaries through observed data
 (判别性模型) model the conditional distribution $p(y|x, \theta)$

Generative Model: include the distribution of data itself
 (生成式模型) model the full joint distribution $p(x, y|\theta)$

Classification Background

target / output: $y \in \{1, 2, \dots, C\}$

input: $X = (X_1, X_2, \dots, X_D)$, where $X_i \in \{1, 2, \dots, K\} \in \mathbb{R}$

$$p(y=c|X_{\text{new}}, \theta) = \frac{p(y=c|\theta) \cdot p(X_{\text{new}}|y=c, \theta)}{\sum_{c=1}^C p(y=c|\theta) p(X_{\text{new}}|y=c, \theta)}$$

class-conditional

joint distribution $p(X_{\text{new}}, y|\theta)$

marginal distribution $p(X_{\text{new}}|\theta)$

$$\Rightarrow \hat{y} = \operatorname{argmax}_c p(y=c|X_{\text{new}}, \theta)$$

joint distribution: $p(x, y|\theta, \pi) = p(y|\pi) \cdot p(x|y, \theta)$

$p(y=c|\pi) = \pi_c, \sum \pi_c = 1$ class-conditional densities.

$p(x|y=c, \theta_c)$ for class $c=1, 2, \dots, C$

Naïve Bayes Classifier (NBC)

To estimate the parameters π_c, θ_c for $c=1, 2, \dots, C$.

Assumption: the features are conditionally independent given the class label

$$p(x|y=c, \theta_c) = \prod_{j=1}^D p(x_j|y=c, \theta_{jc}) \quad \begin{matrix} \text{if not independent,} \\ \text{O}(C2^D) \text{ parameters} \end{matrix}$$

Example: use a Gaussian Model, so $\theta_{jc} = (\mu_{jc}, \sigma_{jc}^2)$

$\sum_{i=1}^K \mu_{jc,i} = 1$ (j : factor $\leq D$, c : choice $\leq C$, i : factor range $\leq K$)
 assume all features binary, $K=1$. $x_j \in \{0, 1\}$. O(CCD) params

Maximum Likelihood for NBC

Data $\{(x_i, y_i)\}_{i=1}^N$ iid. from some joint distribution $p(x, y)$

probability for a single datapoint:

$$p(x_i, y_i|\theta, \pi) = p(y_i|\pi) \cdot p(x_i|\theta, y_i) = \prod_{c=1}^C \pi_c^{\mathbb{I}(y_i=c)} \prod_{c=1}^C \prod_{j=1}^D p(x_{ij}|\theta_{jc})^{\mathbb{I}(y_i=c)}$$

Denote N_c = number of datapoints with $y_i=c$, s.t. $\sum_{c=1}^C N_c = N$
 $\Rightarrow \log p(D|\Theta, \Pi) = \sum_{c=1}^C N_c \log \Pi_c + \sum_{c=1}^C \sum_{j=1}^D \sum_{i:y_i=c} \log p(X_{ij}|\Theta_{jc})$

1. $\max \sum_{c=1}^C N_c \log \Pi_c$ s.t. $\sum_{c=1}^C \Pi_c = 1 \Rightarrow \Pi_c = \frac{N_c}{N}$

2. $\max \sum_{c=1}^C \sum_{j=1}^D \sum_{i:y_i=c} \log p(X_{ij}|\Theta_{jc}) \Rightarrow \widehat{\Theta}_{jc}$ only depends on j^{th} feature with $y_i=c$

Pro: handle missing data at test time, just marginalize it out:

$$p(y=c | X_{\text{new}}, \Theta) = \frac{\Pi_c \cdot \prod_{j=1}^D p(X_{ij} | y=c, \Theta_{jc})}{\sum_{c'=1}^C p(y=c' | \Theta) \prod_{j=1}^D p(X_{ij} | y=c', \Theta_{jc})}$$

Loosen the independence assumption:

$$p(x | y=c, \Theta_c) = N(x | \mu_c, \Sigma_c)$$

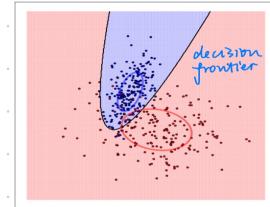
($x \in \mathbb{R}^p$, model class-conditional density as a multivariate normal distribution with Σ_c).

$$\begin{aligned} p(y=c | X_{\text{new}}, \Theta) &= \frac{p(y=c | \Theta) \cdot p(X_{\text{new}} | y=c, \Theta)}{\sum_{c=1}^C p(y=c | \Theta) p(X_{\text{new}} | y=c, \Theta)} \\ &= \frac{\Pi_c / 2\pi |\Sigma_c|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c))}{\sum_{c=1}^C \Pi_c / 2\pi |\Sigma_c|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c))} \end{aligned}$$

the boundary between classes:

$$\frac{\Pi_c / 2\pi |\Sigma_c|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c))}{\Pi_c' / 2\pi |\Sigma_c'|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x - \mu_c')^T \Sigma_c'^{-1} (x - \mu_c'))} = 1$$

— quadratic discriminant analysis (QDA)

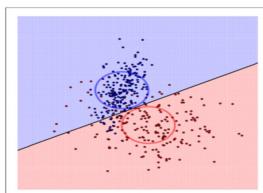


(*) Special case: Shared / tied covariance matrices across different classes

$$\begin{aligned} p(y=c | x, \Theta) &\propto \Pi_c \exp(-\frac{1}{2}(x - \mu_c)^T \Sigma^{-1} (x - \mu_c)) \\ &= \exp(\underbrace{\mu_c^T \Sigma^{-1} x - \frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \log \Pi_c}_{{\beta}_c = \Sigma^{-1} \mu_c \text{ denote as } {\beta}_c} \cdot \exp(-\frac{1}{2} x^T \Sigma^{-1} x)) \end{aligned}$$

$$\begin{aligned} &= \exp({\beta}_c^T x + t_c) \Rightarrow \frac{\exp({\beta}_c^T x + t_c)}{\sum_c \exp({\beta}_c^T x + t_c)} = \text{softmax}(\eta)_c \\ &, \eta = [{\beta}_1^T x + t_1, \dots, {\beta}_C^T x + t_C] \end{aligned}$$

— linear discriminant analysis (LDA)



softmax maps a set of numbers to a probability distribution with mode at the maximum, $\text{softmax}([a, b, c]) = ([p_1, p_2, p_3])$

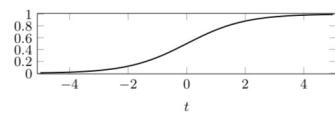
(*) Two class LDA

when we have only 2 classes, say 0 and 1.

$$\begin{aligned} P(y=1|X, \theta) &= \frac{\exp(\beta_1^T X + \gamma_1)}{\exp(\beta_1^T X + \gamma_1) + \exp(\beta_0^T X + \gamma_0)} \\ &= \frac{1}{1 + \exp[-((\beta_1 - \beta_0)^T X + (\gamma_1 - \gamma_0))]} \\ &= \text{Sigmoid } ((\beta_1 - \beta_0)^T X + (\gamma_1 - \gamma_0)) \end{aligned}$$

sigmoid function:

$$\text{sigmoid}(t) = \frac{1}{1 + e^{-t}}$$



MLE for QDA and LDA:

$$\sum \pi_c = \frac{N_c}{N}$$

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} X_i$$

$$\hat{\Sigma}_c = \frac{1}{N_c} \sum_{i:y_i=c} (X_i - \hat{\mu}_c)(X_i - \hat{\mu}_c)^T$$

(number of parameters: $O(CD^2)$).

Logistic Regression

Binary classification: 0 and 1

This is a discriminative method: $p(y|w, x)$

$Z \sim \text{Bernoulli}(\theta)$, $\theta \in [0, 1]$: $\begin{cases} \zeta = 1 & \text{with prob} = \theta \\ \zeta = 0 & \text{with prob} = 1 - \theta \end{cases}$, $p(x|\theta) = \theta^x (1-\theta)^{1-x}$

model: $y \sim \text{Bernoulli}(f(x, w))$, where $f \in [0, 1]$, x is input, w is param
logistic regression builds up on a linear model, composed with sigmoid:

$$p(y|w, x) = \text{Bernoulli}(y| \text{sigmoid}(w \cdot x))$$

$$\left\{ \begin{array}{l} p(y_{\text{new}}=1 | X_{\text{new}}, w) = \text{sigmoid}(w \cdot X_{\text{new}}) \\ \hat{y}_{\text{new}} = \mathbb{I}(\text{sigmoid}(w \cdot X_{\text{new}}) \geq \frac{1}{2}) = \mathbb{I}(w \cdot X_{\text{new}} \geq 0) \end{array} \right.$$

$$p(y | X, w) = \prod_{i=1}^N \tau(w^T x_i)^{y_i} \cdot (1 - \tau(w^T x_i))^{1-y_i}, y_i \in \{0, 1\}$$

$$\Rightarrow \text{NLL}(y | X, w) = -\sum_{i=1}^N (y_i \log \mu_i + (1-y_i) \log(1-\mu_i)), \mu_i = \tau(w^T x_i)$$

the cross-entropy between y_i and μ_i

Entropy $H(X) = -\sum_x p(x) \log p(x)$, Bernoulli has $H(x) = -\theta \log \theta - (1-\theta) \log(1-\theta)$

$$\nabla_w \text{NLL}(y | X, w) = \sum_{i=1}^N x_i (\mu_i - y_i) = X^T (\mu - y)$$

$\left\{ \begin{array}{l} \text{Hessian: } H = X^T S X, S \text{ is a diagonal matrix where } S_{ii} = \mu_i(1-\mu_i) \\ \text{Gradient: } g = -X^T (\mu - y) \end{array} \right.$

The Newton Update Rule:

$$\begin{aligned} w_{t+1} &= w_t - H_t^{-1} g_t \\ &= w_t + (X^T S_t X)^{-1} X^T (y - \mu) \\ &= (X^T S_t X)^{-1} X^T S_t (X w_t + S_t^{-1} (y - \mu)) \\ &= (X^T S_t X)^{-1} X^T S_t z_t, \text{ where } z_t = X w_t + S_t^{-1} (y - \mu) \end{aligned}$$

Then w_{t+1} is a solution of the following:

$\min \sum_{i=1}^N S_{t,ii} (z_{t,i} - w^T x_i)^2$, weighted least squares

Multiclass Logistic Regression

(also a discriminative classifier)

Input $\mathbf{x} \in \mathbb{R}^D$, output $y \in \{1, 2, \dots, C\}$

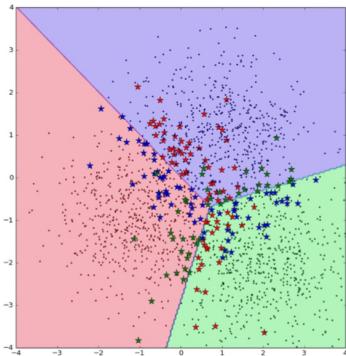
parameters $\mathbf{w}_c \in \mathbb{R}^D$ for every class $c = 1, \dots, C$, denoted as $\mathbf{W}_{D \times C}$

$$p(y=c|\mathbf{x}, \mathbf{W}) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\mathbf{w}_c^\top \mathbf{x})}$$

$$\mathbf{W} = \begin{bmatrix} | & | & | \\ \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_C \\ | & | & | \end{bmatrix}_{D \times C}$$

$$\Rightarrow p(y=c|\mathbf{x}, \mathbf{W}) = \text{softmax}([\mathbf{w}_1^\top \mathbf{x}, \dots, \mathbf{w}_C^\top \mathbf{x}])$$

$$\text{softmax}([\alpha_1, \dots, \alpha_C]^\top) = \left[\frac{e^{\alpha_1}}{Z}, \dots, \frac{e^{\alpha_C}}{Z} \right]^\top, Z = \sum_{c=1}^C e^{\alpha_c}$$



Extension to multiclass:
Sigmoid \rightarrow softmax

Exercises

1. Maximum Likelihood for Logistic Regression

Consider the sigmoid function, defined as $\sigma(z) = \frac{1}{1+e^{-z}}$. Note that $\lim_{z \rightarrow -\infty} \sigma(z) = 0$ and $\lim_{z \rightarrow \infty} \sigma(z) = 1$. Thus, for binary classification problems, we can compose a linear function with the sigmoid function to model the probability that a given input \mathbf{x} belongs to one of the two classes $\{0, 1\}$. More precisely, for parameter vector $\mathbf{w} \in \mathbb{R}^D$, and input vector $\mathbf{x} \in \mathbb{R}^{D, 1}$, the label $y \in \{0, 1\}$ is given by the following model:

$$\Pr(y = 1 \mid \mathbf{w}, \mathbf{x}) = \sigma(\mathbf{x}^\top \mathbf{w}) \quad (2.1)$$

$$\Pr(y = 0 \mid \mathbf{w}, \mathbf{x}) = 1 - \sigma(\mathbf{x}^\top \mathbf{w}) \quad (2.2)$$

$$\nabla'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} (1 - \frac{1}{1+e^{-z}}) = \nabla(z)(1 - \nabla(z))$$

$$P(y_1, y_2, \dots, y_N \mid \mathbf{w}, \mathbf{x}) = \prod_{i:y_i=1} \nabla(\mathbf{x}_i^\top \mathbf{w}) \prod_{i:y_i=0} (1 - \nabla(\mathbf{x}_i^\top \mathbf{w}))$$

$$NLL(y \mid \mathbf{w}, \mathbf{x}) = -\log P(y \mid \mathbf{w}, \mathbf{x}) = -\sum_{i:y_i=1} \log(\nabla(\mathbf{x}_i^\top \mathbf{w})) - \sum_{i:y_i=0} \log(1 - \nabla(\mathbf{x}_i^\top \mathbf{w}))$$

$$\begin{aligned} \nabla_{\mathbf{w}} NLL(y \mid \mathbf{w}, \mathbf{x}) &= -\sum_{i:y_i=1} \frac{\nabla(\mathbf{x}_i^\top \mathbf{w})(1 - \nabla(\mathbf{x}_i^\top \mathbf{w}))}{\nabla(\mathbf{x}_i^\top \mathbf{w})} \mathbf{x}_i + \sum_{i:y_i=0} \frac{\nabla(\mathbf{x}_i^\top \mathbf{w})(1 - \nabla(\mathbf{x}_i^\top \mathbf{w}))}{1 - \nabla(\mathbf{x}_i^\top \mathbf{w})} \mathbf{x}_i \\ &= -\sum_{i:y_i=1} (1 - \nabla(\mathbf{x}_i^\top \mathbf{w})) \mathbf{x}_i + \sum_{i:y_i=0} \nabla(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i \\ &= -\sum_{i:y_i=1} \mathbf{x}_i + \sum_{i=1}^N \nabla(\mathbf{x}_i^\top \mathbf{w}) \mathbf{x}_i \end{aligned}$$

$$\mathbf{H} = \nabla_{\mathbf{w}}^2 NLL(y \mid \mathbf{w}, \mathbf{x}) = \left(\sum_{k=1}^N \nabla(\mathbf{x}_k^\top \mathbf{w})(1 - \nabla(\mathbf{x}_k^\top \mathbf{w})) \mathbf{x}_{ki} \mathbf{x}_{kj} \right)_{i,j}$$

$$\begin{aligned} &= \begin{bmatrix} \mathbf{x}_{11} & \cdots & \mathbf{x}_{1D} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{ND} & \cdots & \mathbf{x}_{ND} \end{bmatrix} \begin{bmatrix} \nabla(\mathbf{x}_1^\top \mathbf{w})(1 - \nabla(\mathbf{x}_1^\top \mathbf{w})) \\ \ddots \\ \nabla(\mathbf{x}_N^\top \mathbf{w})(1 - \nabla(\mathbf{x}_N^\top \mathbf{w})) \end{bmatrix} \begin{bmatrix} \mathbf{x}_{11} & \cdots & \mathbf{x}_{1D} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{ND} & \cdots & \mathbf{x}_{ND} \end{bmatrix} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X} \end{aligned}$$

where \mathbf{S} is a diagonal matrix s.t. $S_{ii} = \nabla(\mathbf{x}_i^\top \mathbf{w})(1 - \nabla(\mathbf{x}_i^\top \mathbf{w}))$

\mathbf{S} is positive definite because

$$\forall \mathbf{y} \in \mathbb{R}^{N \times 1}, \mathbf{y}^\top \mathbf{S} \mathbf{y} = \sum_{i=1}^N y_i^2 S_{ii} \geq 0, \text{ and } \mathbf{y}^\top \mathbf{S} \mathbf{y} = 0 \text{ iff. } \mathbf{y} = 0$$

Thus, \mathbf{H} is positive semi-definite because

$$\forall \mathbf{y} \in \mathbb{R}^{D \times 1}, \mathbf{y}^\top (\mathbf{X}^\top \mathbf{S} \mathbf{X}) \mathbf{y} = (\mathbf{X} \mathbf{y})^\top \mathbf{S} (\mathbf{X} \mathbf{y}) \geq 0.$$

and $(\mathbf{X} \mathbf{y})^\top \mathbf{S} (\mathbf{X} \mathbf{y}) = 0$ iff. $\mathbf{X} \mathbf{y} = 0$, which doesn't necessarily mean $\mathbf{y} = 0$.

Since the Hessian is positive semi-definite, NLL is a convex function of \mathbf{w} , and the obtained solution is the global minimum.

SVM (Support Vector Machines)

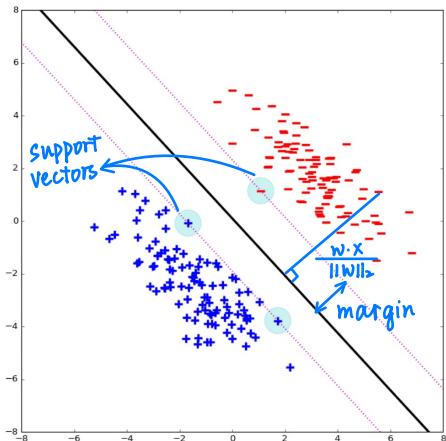
Binary Classification

1. Goal: find a linear separator

data is linearly separable iff. all points are classified correctly

2. Maximum Margin Principle: max the distance of the closest point from the decision boundary

support vectors



$$\text{Hyperplane: } H \equiv w \cdot x + w_0 = 0$$

$$\text{Distance between } x \in \mathbb{R}^D \text{ and } H: \frac{|w \cdot x + w_0|}{\|w\|_2}$$

Separable Case:

$$\min \frac{1}{2} \|w\|^2 \text{ s.t. } y_i(w \cdot x_i + w_0) \geq 1 \quad (\text{where } y_i \in \{-1, 1\})$$

quadratic program

Non-separable Case:

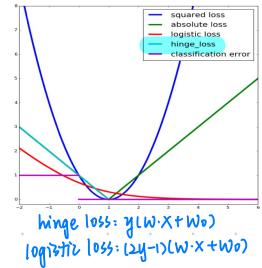
$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \quad \begin{array}{l} \text{(regularizer)} \\ \text{(loss function)} \end{array}$$

$\xi_i = \max\{0, 1 - y_i(w \cdot x_i + w_0)\}$ hinge loss function

s.t. $y_i(w \cdot x_i + w_0) \geq 1 - \xi_i$ ← slack to violate constraints

Logistic Regression: Loss Function

$$\text{NLL}(y_i; w, x_i) = -y_i \log\left(\frac{1}{1+e^{-w \cdot x}}\right) + (1-y_i) \log\left(\frac{1}{1+e^{-(w \cdot x)}}\right)$$
$$= \log\left(1+e^{-z_i(w \cdot x_i)}\right) = \log\left(1+e^{-(2y_i-1)(w \cdot x_i)}\right)$$



Multiclass Classification

- one-vs-one: train $\binom{K}{2}$ classifiers, choose the most commonly occurring y
- one-vs-rest: train K classifiers (one class vs the rest)

One-vs-One

- Training roughly $K^2/2$ classifiers
- Each training procedure only uses on average $2/K$ of the training data
- Resulting learning problems are more likely to be "natural"

One-vs-Rest

- Training only K classifiers
- Each training procedure uses all the training data
- Resulting learning problems are less likely to be "natural"

Exercises

1. SVM problems:

$$(a) \text{minimize } \frac{1}{2} \|w\|_2^2$$

$$\text{subject to: } y_i(w \cdot x_i + w_0) \geq 1$$

$$(b) \text{maximize } \alpha$$

$$\text{subject to: } y_i(w \cdot x_i + w_0) \geq \alpha, \|w\|_2 = 1$$

$$(a) \Leftrightarrow (b)$$

Suppose you use the primal SVM formulation for the non-separable case, i.e., with slack variables ζ_i , but your data is actually linearly separable. Do you always recover the "true" max-margin separating hyperplane?

No, not necessarily.

In non-separable case:

$$\text{minimum } \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \zeta_i$$

$$\text{subject to: } y_i(w \cdot x_i + w_0) \geq 1 - \zeta_i$$

Correctly classified examples: $0 \leq \zeta_i \leq 1$

Misclassified examples: $\zeta_i > 1$

If we set the penalty C to a very small number, the motivation to minimize $\|w\|_2$ weighs more than that to separate the data correctly, which may not be the "true" hyperplane.

2. Kernel methods

Kernel methods are widely applicable in machine learning and not just restricted to support vector machines. For instance, linear regression using kernel basis expansion, is nothing but a *kernelised* form of linear regression. Let us briefly revisit the Ridge regression objective:

$$\mathcal{L}_{\text{ridge}}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

We observe that the optimal solution \mathbf{w} can always be written as $\mathbf{w} = \sum_{i=1}^N \alpha_i \mathbf{x}_i$. Why? Suppose not. Let us write $\mathbf{w} = \mathbf{w}_s + \mathbf{w}_\perp$, where \mathbf{w}_s is the component of \mathbf{w} that lies in the span $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and \mathbf{w}_\perp is orthogonal to this linear subspace spanned by the inputs. Then, since $\mathbf{w}_\perp \cdot \mathbf{x}_i = 0$ for all $i = 1, \dots, N$, \mathbf{w}_\perp cannot affect $(\mathbf{w} \cdot \mathbf{x}_i - y_i)^2$ at all. However, $\|\mathbf{w}\|_2^2 = \|\mathbf{w}_s\|_2^2 + \|\mathbf{w}_\perp\|_2^2$. Thus, setting $\mathbf{w}_\perp = \mathbf{0}$ can only decrease $\mathcal{L}_{\text{ridge}}$. Thus, equivalently we can try to solve for α_i rather than \mathbf{w} , and we get a *kernelised* form of linear regression. You will work out the details on Problem Sheet 3.

Dual Formulation of SVM

primal form : minimize $F(\mathbf{z}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i$ here concave
 subject to $g_i(\mathbf{z}) \geq 0 \quad (i=1, 2, \dots, m) \quad \text{convex}$
 $h_j(\mathbf{z}) = 0 \quad (j=1, 2, \dots, l) \quad \text{affine} \quad (\mathbf{w}_0 + \mathbf{w}_1 \mathbf{z}_1 + \dots + \mathbf{w}_l \mathbf{z}_l)$

Lagrange Function : $\Lambda(\mathbf{z}; \alpha, \mu) = F(\mathbf{z}) - \sum_{i=1}^m \alpha_i g_i(\mathbf{z}) - \sum_{j=1}^l \mu_j h_j(\mathbf{z})$

For convex problems, Karush-Kuhn-Tucker (KKT) conditions provide necessary and sufficient conditions for a critical point of Λ to be the minimum of the original constrained optimization problem $(\mathbf{z}^*, \alpha^*, \mu^*)$

1. Primal Feasibility : $g_i(\mathbf{z}^*) \geq 0, h_j(\mathbf{z}^*) = 0, i=1, 2, \dots, m; j=1, 2, \dots, l \quad * \nabla_{\mathbf{w}, b} \Lambda(\mathbf{w}, b, \alpha) = 0$
2. Dual Feasibility : $\alpha_i^* \geq 0, i=1, 2, \dots, m$
3. Complementary Slackness : $\alpha_i^* g_i(\mathbf{z}^*) = 0, i=1, 2, \dots, m$

$$\Rightarrow \Lambda(\mathbf{w}, \mathbf{w}_0, \xi; \alpha, \mu) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + \mathbf{w}_0) - 1 - \xi_i) - \sum_{i=1}^N \mu_i \xi_i$$

$$\left\{ \begin{array}{l} \frac{\partial \Lambda}{\partial \mathbf{w}_0} = - \sum_{i=1}^N \alpha_i y_i \\ \frac{\partial \Lambda}{\partial \xi_i} = C - \alpha_i - \mu_i \end{array} \right. \Rightarrow FOC = 0, \quad 0 \leq \alpha_i \leq C$$

$$\nabla_{\mathbf{w}} \Lambda = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\text{Thus } \Rightarrow g(\alpha) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \cdot \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right)$$

dual form : maximize $\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$ concave, N variables
 subject to $\sum_{i=1}^N \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C \quad (N^2 \text{ terms to express})$

Solution : $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad \alpha_i = 0 \text{ or } y_i (\mathbf{w} \cdot \mathbf{x}_i + \mathbf{w}_0) = 1 - \xi_i$
 $\mathbf{w}_0 = y_j - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_j \quad (\text{choose one } \alpha_j \neq 0)$

Prediction : $\mathbf{w} \cdot \mathbf{x}_{\text{new}} = \sum_{i=1}^N \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}_{\text{new}}) \Rightarrow \text{thus Kernel } K(\cdot, \cdot) \text{ is sufficient}$

time complexity = $O(d \log d) = O(d \log d)$

1. Polynomial Basis Expansion : $K(x, x') = \phi_d(x) \cdot \phi_d(x') = (1 + x \cdot x')^d$
 for $d=2$, $\phi_2([x_1, x_2]^T) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2]^T \quad O(d^d) \text{ variables}$
 $(d: \text{input features}, d: \text{power})$

2. Mercer Kernels ($K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, for inputs $\langle \mathbf{x}_i \rangle_{i=1}^N$, each $\mathbf{x}_i \in \mathcal{X}$)

Gram matrix $K = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \dots & K(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$

$g(\alpha) = \mathbf{1}^T \alpha - \frac{1}{2} (\alpha \otimes \alpha)^T K (\alpha \otimes \alpha)$
 $\Rightarrow \text{concave if } K \text{ is positive definite}$

is always positive semi-definite

(Here, $K(x, x')$ is some measure of similarity between x and x')

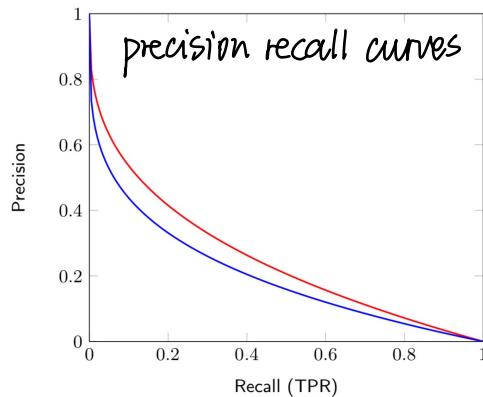
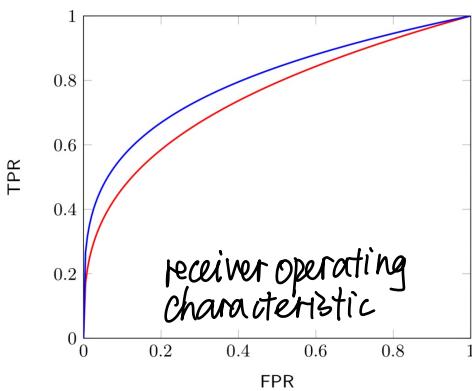
3. Gaussian / RBF Kernels : $K(x, x') = \exp\left(-\frac{\|x-x'\|_2^2}{2\sigma^2}\right) \quad \text{bandwidth}$

Measuring Performance

prediction	(Actual Labels)		Type I errors ↑
	Y	N	
yes	true positive	false positive	⇒ Precision: $P = \frac{TP}{TP+FP}$
no	false negative	true negative	Type II errors
	↓	↓	
	True Positive Rate	False Positive Rate	
	$TPR = \frac{TP}{TP+FN}$	$FPR = \frac{FP}{FP+TN}$	

For multi-class classification, it's common to write confusion matrix

prediction	Actual Labels		K
	1	2	
1	N_{11}	N_{12}	...
2	N_{21}	N_{22}	...
⋮	⋮	⋮	⋮
K	N_{K1}	N_{K2}	...



Neural Network

slides 14~16 - Neural Networks