

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

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

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

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

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

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

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

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

$$R(X^T)$$

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

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

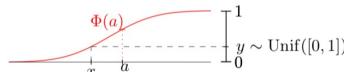
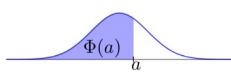
$$= (X^T X)^{-1} X^T \sigma^2 I_n X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

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

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

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

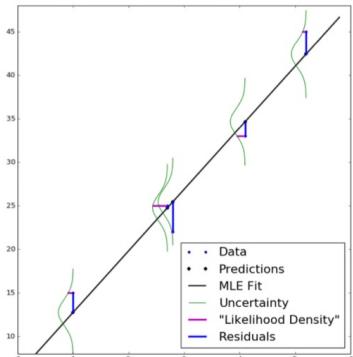
Maximum Likelihood

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

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

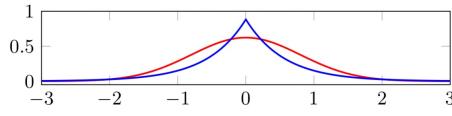
Outliers and Laplace Distribution

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

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

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

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



Fitting, Regularization and Bayesian

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

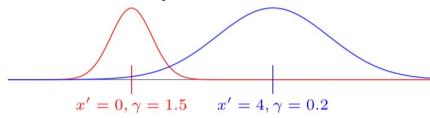
Basis Function Expansion

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

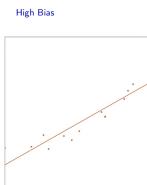
3. model: $y = w_0 + w_1 K(\mu_1, X) + \dots + w_m K(\mu_m, X) + \epsilon = w \cdot \phi(X) + \epsilon$

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)

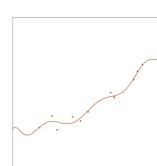
{ overfitting: high variance (training \approx test errors)



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)

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

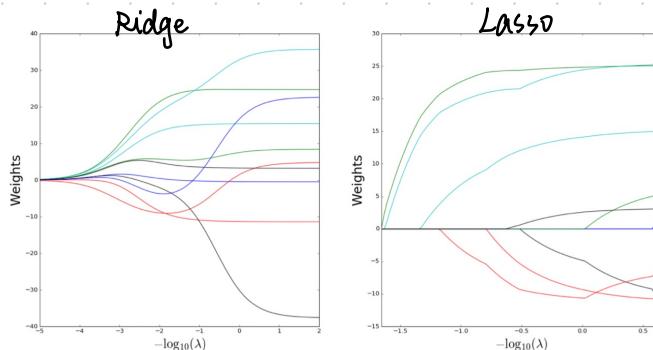
L_2 -regularization
(weight-decay)

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

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:

$\Theta(n^2)$

1. Set set of selected features to $F := \emptyset$

2. 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

3. 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

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)

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

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

prior: $p(\mu_s) = \mathcal{N}(\mu_{\text{a.s.}}, \tau_{\text{a.s.}}^2)$

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

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

Bernstein-von Mises Theorem:

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

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

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

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

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

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

$$\text{posterior: } \underbrace{\text{pdf}(w|D)}_{\text{likelihood}} \propto \underbrace{\mathcal{N}(y|Xw, \Sigma)}_{\text{prior}} = \mathcal{N}(w|w_N, \Sigma_N)$$

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

} complexity!

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

$$= \mathcal{N}(y|w_N^T X_{\text{new}}, \tau^2 + X_{\text{new}}^T \Sigma_N X_{\text{new}}) \quad \begin{matrix} \text{aleatoric} \\ \text{epistemic} \end{matrix}$$

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

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

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

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

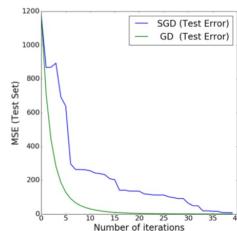
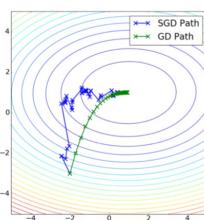
$$w_{t+1} = w_t - H_t^{-1} g_t \quad f_{\text{quad}}(w) = f(w_t) + g_t^T (w - w_t) + \frac{1}{2} (w - w_t)^T H_t (w - w_t)$$
$$\nabla_w f_{\text{quad}} = g_t + H_t (w - w_t) := 0$$

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