

Probability

Bayes Theorem:
P(Y = ±1|X) = P(X|Y=±1)P(Y=±1) / (P(X|Y=+1)P(Y=+1)+P(X|Y=-1)P(Y=-1))

X and Y are independent iff P(X, Y) = P(X)P(Y)
X and Y are uncorrelated iff E(X, Y) = E(X)E(Y)

Matrix calculus

f(x) = Ax + x'A + x'x + x'Ax => df(x)/dx = A' + A + 2x + Ax + A'x
Hi,j = d^2 f / (dxi dxj); grad(ax) = aI; J = |dxi/dyi| <=> J^-1 = |dyi/dxi|

Perceptron

f(x) = theta . x + theta_0 = sum_{i=1}^d theta_i x_i + theta_0, y-hat = { 1, if f(x) >= 0; -1, if f(x) < 0

Decision boundary, a hyperplane in R^d:
H = {x in R^d : f(x) = 0} = {x in R^d : theta . x + theta_0 = 0}

theta is the normal of the hyperplane,
theta_0 is the offset of the hyperplane from origin,
- theta_0 / ||theta|| is the signed distance from the origin to hyperplane.

Perceptron algorithm,
Input: (x1, y1), ..., (xn, yn) in R^d x {±1}
while some yi != sign(theta . xi)
 pick some misclassified (xi, yi)
 theta <- theta + yi xi

Given a linearly separable data, perceptron algorithm will take no more than R^2 / gamma^2 updates to converge, where R = max_i ||xi|| is the radius of the data, gamma = min_i yi (theta . xi) / ||theta|| is the margin.
Also, theta . x / ||theta|| is the signed distance from H to x in the direction theta.

theta = sum_i alpha_i y_i x_i, thus any inner product space will work, this is a kernel.

Gradient descent view of perceptron, minimize margin cost function J(theta) = sum_i (-yi(theta . xi))_+ with theta <- theta - eta grad J(theta)

Support Vector Machine

Hard margin SVM,
min_theta ||theta||^2, such that yi theta . xi >= 1 (i = 1, ..., n)
Soft margin SVM,
min_theta ||theta||^2 + C sum_{i=1}^n (1 - yi theta . xi)_+

Regularization and SVMs: Simulated data with many features phi(x); C controls trade-off between margin 1/||theta|| and fit to data; Large C: focus on fit to data (small margin is ok). More overfitting. Small C: focus on large margin, less tendency to overfit. Overfitting increases with: less data, more features.

theta = sum_j alpha_j y_j x_j, alpha_j != 0 only for support vectors.

Width of the margin is 2 / ||theta||.

K(xi_i, x_j) = phi(xi_i) . phi(x_j), K is called a kernel.
Solve alpha_j to determine sum_j alpha_j y_j phi(x_j).
Compute the classifier for a test point x via
theta . phi(x) = sum_j alpha_j y_j K(x_j, x)

degree-m polynomial kernel: Km(x, x-tilde) = (1 + x . x-tilde)^m
RBF kernel (infinite dimation): Krbf(x, x-tilde) = exp(-gamma ||x - x-tilde||^2)

Decision Theory

Loss function: l : Y x Y -> R, and l(y-hat, y) is the cost of predicting y-hat when the outcome is y.

Risk for a given class: R(alpha_i|x) = sum_{i=1}^C lambda_ij P(w = j|x)

Assume (X, Y) are chosen i.i.d according to some probability distribution on X x Y. Risk is misclassification probability:
R(f) = El(f(X), Y) = Pr(f(X) != Y)

Bayes Decision Rule is
f*(x) = { 1, if P(Y = 1|x) > P(Y = -1|x); -1, otherwise.
and the optimal risk (Bayes risk) R* = inf_f R(f) = R(f*)

Excess risk is for any f : X -> {-1, +1},
R(f) - R* = E(1[f(x) != f*(x)]|2P(Y = +1|X) - 1|)

Risk in Regression is expected squared error:
R(f) = El(f(X), Y) = EE[f(X) - Y^2|X]

Bias-variance decomposition:
R(f) = E[(f(X) - E[Y|X])^2] + E[(E[Y|X] - Y)^2]
= E[(f(X) - f*(X))^2] + E[(f*(X) - Y)^2]
R(f) = E[(f(X) - f*(X))^2] + R(f*)
R(f) - R(f*) = E[(f(X) - f*(X))^2], f*(X) = E[Y|X]

Generative and Discriminative

Discriminative models: P(X, Y) = P(X)P(Y|X).
Estimate P(Y|X), then pretend out estimate P-hat(Y|X) is the actual P(Y|X) and plug in bayes rule expression.

Generative model: P(X, Y) = P(Y)P(X|Y).
Estimate P(Y) and P(X|Y), then use bayes theorem to calculate P(Y|X) and use discriminative model.

Gaussian class conditional densities P(X|Y = +1), P(X|Y = -1) (with the same variance), the posterior probability is logistic:
P(Y = +1|x) = 1 / (1 + exp(-x . beta - beta_0))
beta = Sigma^-1 (mu_1 - mu_0), beta_0 = (mu_0' Sigma^-1 mu_0 - mu_1' Sigma^-1 mu_1) / 2 + log (P(Y=1) / P(Y=0))

Estimation

Method of moments: Match moments of the distribution to momemnts measured in the data.

Maximum likelihood: Choose parameter so that the distribution it defines gives the observed data the highest probability (likelihood).

Maximum log likelihood: Log of maximum likelihood, equivalent to maximum likelihood since log is monotonically increase; it is useful since it can change product to sum.

Penalized maximum likelihood: Add a penalty term in the maximum (log) likelihood equation; treat the penalty term as some imaginary data points crafted for desired probability.

Bayesian estimate: Treat parameter as a random variable, then update based on observed value (data).
Prior: pi(p) = 1,
Posterior: P(p|X1 = 1) = P(X1 = 1|p)pi(p) / integral P(X1 = 1|q)dpi(q)

Maximum a posterior probability: the mode of the posterior.
If uniform prior, MAP is MLE; if not uniform prior, MAP is Penalized MLE.

Multivariate Normal Distribution

x in R^d : p(x) = 1 / ((2pi)^d / 2 |Sigma|^(1/2)) * exp(-1/2 (x - mu)' Sigma^-1 (x - mu))

Covariance matrix: Sigma = E(X - mu)(X - mu)'
Symmetric: Sigma_i,j = Sigma_j,i
Non-negative diagonal entries: Sigma_i,i >= 0
Positive semidefinite: for all v in R^d, v' Sigma v >= 0

Super-level sets of pdf:
xi_r = {x in R^d : (x - mu)' Sigma^-1 (x - mu) <= r^2}.
Volume of xi_r proportional to product_{i=1}^d sigma_i = sqrt(|Sigma|)

Spectral Theorem for non-diagonal covariance:
U = [v1, v2, ..., vn], Lambda = diag([lambda1, lambda2, ..., lambda_n]')
We can eigen decompose Sigma^-1 = U Lambda^-1 U', this is like to change to a different eigen spaces, where covariances (Lambda) diagonal axis-alianed.

Assume independent,
N(mu_x, Sigma) + N(mu_y, Sigma_y) = N(mu_x + mu_y, Sigma_x + Sigma_y)

Given a d-dimensaional Gaussian X ~ N(mu, Sigma),
write X = [Y; Z], mu = [mu_Y; mu_Z], Sigma = [Sigma_YY Sigma_YZ; Sigma_ZY Sigma_ZZ],
where Y in R^m, and Z in R^(d-m). Then Y ~ N(mu_Y, Sigma_YY)

Given a d-dimensaional Gaussian X ~ N(mu, Sigma),
matrix A in R^m x d and vector b in R^m, define Y = AX + b.
Then Y ~ N(A mu + b, A Sigma A')

Given a d-dimensaional Gaussian X ~ N(mu, Sigma),
with Sigma positive definite,
Y = Sigma^-1/2 (X - mu) ~ N(0, I)

Gaussian maximum likelihood estimation:
Sample mean: mu-hat = 1/n sum_{i=1}^n xi_i;
Sample covariance: Sigma-hat = 1/n sum_{i=1}^n (xi_i - mu-hat)(xi_i - mu-hat)'

Linear Regression

Given $\mathbf{X} \in \mathbb{R}^p$, $Y \in \mathbb{R}$, consider linear(affine) prediction rules, $F_{\text{lin}} := \{\mathbf{x} \mapsto \mathbf{x}'\boldsymbol{\beta} + \beta_0 : \boldsymbol{\beta} \in \mathbb{R}^p, \beta_0 \in \mathbb{R}\}$

Empirical risk minimization

Empirical risk is the sample average of squared error:
 $\hat{R}(f) = \hat{\mathbb{E}}_n \ell(f(\mathbf{X}), Y) = \frac{1}{n} \sum_{i=1}^n n(f(\mathbf{X}_i) - Y_i)^2$
Choose $\hat{f} := \arg \min_{f \in F_{\text{lin}}} \hat{\mathbb{E}}_n \ell(f(\mathbf{X}), Y)$

Find $\hat{f} : \mathbf{x} \mapsto \mathbf{x}'\hat{\boldsymbol{\beta}}$, such that
 $\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \underbrace{\sum_{i=1}^n (\mathbf{X}'_i \boldsymbol{\beta} - Y_i)^2}_{\text{RSS}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \underbrace{\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|^2}_{\text{RSS}}$

where **design matrix** $\mathbf{X} \in \mathbb{R}^{n \times p}$ and **response vector** $\mathbf{y} \in \mathbb{R}^n$.

Normal equations: $\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\mathbf{y} \Rightarrow \hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$

Projection Theorem also leads to normal equations:
 $(\mathbf{y} - \hat{\mathbf{y}})^{-1}\mathbf{X} = 0 \Leftrightarrow \mathbf{X}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0 \Leftrightarrow \mathbf{X}'\mathbf{y} = \mathbf{X}'\mathbf{X}\boldsymbol{\beta}$

Linear model with additive Gaussian noise

Model the conditional distribution of Y given $\mathbf{X} = \mathbf{x}$ as:
 $P(Y|\mathbf{X} = \mathbf{x}) = \mathcal{N}(\mathbf{x}'\boldsymbol{\beta}, \sigma^2)$
Equivalently: $Y = \mathbf{x}'\boldsymbol{\beta} + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2)$

Maximum likelihood is least square:
 $L(\boldsymbol{\beta}) = \prod_{i=1}^n p(Y_i|\mathbf{X}_i, \boldsymbol{\beta}) \Leftrightarrow \ell(\boldsymbol{\beta}) = g(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{X}'_i \boldsymbol{\beta})^2$

Fix \mathbf{X} . Provided $\mathbb{E}\mathbf{y} = \mathbf{X}\boldsymbol{\beta}$ and $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{I}$

Bayesian analysis: Treat $\boldsymbol{\beta}$ as a r.v. with prior distribution $\mathcal{N}(\mathbf{0}, \tau^2\mathbf{I})$, then compute posterior distribution $P(\boldsymbol{\beta}|\mathbf{X}, Y)$.

$P(\boldsymbol{\beta}|\mathbf{X}_1, Y_1, \dots, \mathbf{X}_n, Y_n) \propto P(Y_1, \dots, Y_n|\boldsymbol{\beta}, \mathbf{X}_1, \dots, \mathbf{X}_n)P(\boldsymbol{\beta})$
 $P(\boldsymbol{\beta}|\mathbf{X}_1, Y_1, \dots, \mathbf{X}_n, Y_n) \propto \exp(-\frac{1}{2}(\sum_{i=1}^n \frac{(Y_i - \mathbf{X}'_i \boldsymbol{\beta})^2}{\sigma^2} + \frac{1}{\tau^2}\|\boldsymbol{\beta}\|^2))$

Linear Regression Regularization

Trading off bias and variance: some increase in bias can give a big decrease in variance.

Subset selection is like $L0$ regularization: RSS decreases as the complexity increases because the best fit with a smaller subset is always possible with a larger subset.

Find a path through subset space: using cross-validation and forward-stepwise selection or backward-stepwise selection (need $n > p$).

Ridge regression is like $L2$ regularization:
 $\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} (\sum_{i=1}^n (y_i - \mathbf{x}'_i \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p p\beta_j^2)$
 $\hat{\boldsymbol{\beta}}^{\text{ridge}} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{x}'\mathbf{y}$

Lasso regression is like $L1$ regularization:
 $\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} (\sum_{i=1}^n (y_i - \mathbf{x}'_i \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p p|\beta_j|)$

While ridge regression leads to reduced, but non-zero values of the coefficients, Lasso regression forces some coefficients to be zero.

Bayesian analysis: Ridge regression is equivalent to a MAP estimate with a gaussian prior. Lasso regression is equivalent to a MAP estimate with a Laplace prior.

Logistic Regression

Model **log odds** $(\log p/(1 - p))$ as an affine function of \mathbf{x} .

$P(Y = 1|\mathbf{x}) = \frac{1}{1 + \exp(\boldsymbol{\beta}'\mathbf{x})}$ Given data $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n) \in \mathbb{R}^p \times \{0, 1\}$, estimate $\boldsymbol{\beta}$ with maximum likelihood.

Log likelihood:
 $\ell(\boldsymbol{\beta}) = \sum_{i=1}^n y_i \log \mu_i(\boldsymbol{\beta}) + (1 - y_i) \log(1 - \mu_i(\boldsymbol{\beta}))$,
where $\mu_i(\boldsymbol{\beta}) = P(Y = 1|\mathbf{X} = \mathbf{x}_i, \boldsymbol{\beta}) = \frac{1}{1 + \exp(-\boldsymbol{\beta}'\mathbf{x}_i)}$

$\nabla_{\boldsymbol{\beta}} \mu_i(\boldsymbol{\beta}) = \mu_i(\boldsymbol{\beta})(1 - \mu_i(\boldsymbol{\beta}))\mathbf{x}_i$
 $\nabla_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \mu_i(\boldsymbol{\beta}))\mathbf{x}_i = \mathbf{X}'(\mathbf{y} - \boldsymbol{\mu})$
 $\nabla_{\boldsymbol{\beta}}^2 \ell(\boldsymbol{\beta}) = \sum_{i=1}^n -\mu_i(\boldsymbol{\beta})(1 - \mu_i(\boldsymbol{\beta}))\mathbf{x}_i \mathbf{x}'_i = -\mathbf{X}'\text{diag}(\boldsymbol{\mu}(1 - \boldsymbol{\mu}))\mathbf{X}$
 $\hat{\boldsymbol{\beta}}^{\text{ml}}$ solves: $\sum_{i=1}^n y_i \mathbf{x}_i = \sum_{i=1}^n \mu_i \boldsymbol{\beta} \mathbf{x}_i$

Gradient ascent:
 $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + \eta \nabla_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}^{(t)}) : O(np)/\text{step}$
Stochastic gradient ascent:
 $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + \eta (y_{i_t} - \mu_{i_t}(\boldsymbol{\beta}^{(t)}))\mathbf{x}_{i_t} : O(p)/\text{step}$
Newton-Raphson method:
 $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - [\nabla_{\boldsymbol{\beta}}^2 \ell(\boldsymbol{\beta}^{(t)})]^{-1} \nabla_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}^{(t)})$

Newton's method for root finding: $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$

Prediction $\hat{p}(y|\mathbf{x}) = \begin{cases} P(Y = 1|\mathbf{x}), & \text{if } y = 1 \\ P(Y = -1|\mathbf{x}), & \text{if } y = -1 \end{cases}$

Log loss (Binomial Deviance): $\ell_{\log}(\hat{p}(\cdot|\mathbf{x}), y) = -\log(\hat{p}(y|\mathbf{x}))$
Minimize: $\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \boldsymbol{\beta}'\mathbf{x}_i))$

Linear Discriminant Analysis

Linear discriminant functions:
 $\delta_k(\mathbf{x}) = \boldsymbol{\mu}'_k \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}'_k \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \log \pi_k$

Estimate with Maximum likelihood:
 $\pi_k = P(Y = k) \Leftrightarrow \hat{\pi}_k = \frac{n_k}{n}$
 $\boldsymbol{\mu}_k = \mathbb{E}[\mathbf{X}|Y = k] \Leftrightarrow \hat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i: y_i = k} \mathbf{x}_i$
 $\boldsymbol{\Sigma} = \text{Var}[\mathbf{X}|Y = k] \Leftrightarrow \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_k \sum_i i : y_i = k (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)'$

SVM with Convex Optimization

Lagrangian: rewrite constraint as penalties for a convex optimization problem such that $L(x, \lambda) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x)$.

Weak duality: $p^* = \min_x \max_{\lambda \geq 0} L(x, \lambda) \geq \max_{\lambda \geq 0} \min_x L(x, \lambda) = d^*$

$\underbrace{\hspace{10em}}$
primal

$\underbrace{\hspace{10em}}$
dual

Strong duality:
if there is a saddle point (x^*, λ^*) such that for all x and $\lambda \geq 0$, $L(x^*, \lambda) \leq L(x^*, \lambda^*) \leq L(x, \lambda^*)$, then primal and dual have the same value ($p^* = d^*$).

Karush-Kuhn-Tucker optimality conditions:
Primal feasibility: $f_i(x) \leq 0$; Dual feasibility: $\lambda_i \geq 0$
Complementary slackness: $\lambda_i f_i(x) = 0$
Stationarity: $\nabla f_0(x) + \sum_i \lambda_i \nabla f_i(x) = 0$

Hard margin SVM:
 $L(\boldsymbol{\theta}, \alpha) = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i \boldsymbol{\theta}'\mathbf{x}_i)$
 $g(\alpha) = \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \alpha)$
setting $\boldsymbol{\theta}^* = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$,
 $g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}'_i \mathbf{x}_j$

Hard margin SVM dual problem:
 $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}'_i \mathbf{x}_j$, s.t. $\alpha_i \geq 0$ ($i = 1, \dots, n$).
 $\min_{\alpha} \frac{1}{2} \alpha' \text{diag}(\mathbf{y}) \mathbf{K} \text{diag}(\mathbf{y}) \alpha - \alpha' \mathbf{1}$, s.t. $\alpha \geq \mathbf{0}$.

Soft margin SVM:
 $L(\boldsymbol{\theta}, \xi, \alpha, \lambda) = \frac{1}{2} \|\boldsymbol{\theta}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i \boldsymbol{\theta}'\mathbf{x}_i - \xi_i) - \sum_{i=1}^n \lambda_i \xi_i$

Soft margin SVM dual problem:
 $\min_{\alpha} \frac{1}{2} \alpha' \text{diag}(\mathbf{y}) \mathbf{K} \text{diag}(\mathbf{y}) \alpha - \alpha' \mathbf{1}$, s.t. $\mathbf{0} \leq \alpha \leq \frac{C}{n}$.

K-Nearest Neighbors

Given \mathbf{x}_q , take vote among its k nearest neighbors, or take mean of f values of k nearest neighbors if real-values.
 $\hat{f}(\mathbf{x}_q) \leftarrow \frac{1}{k} \sum_{i=1}^k f(\mathbf{x}_i)$

Distance metrics:
p-norm: $\|\mathbf{z}\|_p = (\sum_{i=1}^d |z_i|^p)^{1/p}$