

Math Review

X and **Y** are **independent** iff $\mathbb{P}(\mathbf{X}, \mathbf{Y}) = \mathbb{P}(\mathbf{X}) \mathbb{P}(\mathbf{Y})$
X and **Y** are **uncorrelated** iff $\mathbb{E}(\mathbf{X}, \mathbf{Y}) = \mathbb{E}(\mathbf{X}) \mathbb{E}(\mathbf{Y})$
Expected value of $g(X)$: $E[g(X)] = \int_{-\infty}^{\infty} g(x) f(x) dx$
Variance $\sigma^2 = E[(X - \mu)^2] = E[X^2] - \mu^2$
Determinant of matrix is product of its eigenvalues.

$f(\mathbf{w}) = \mathbf{A}\tilde{\mathbf{x}} + \tilde{\mathbf{x}}^\top \mathbf{A} + \tilde{\mathbf{x}}^\top \tilde{\mathbf{x}} + \tilde{\mathbf{x}}^\top \mathbf{A}\tilde{\mathbf{x}} \Rightarrow \frac{df(\tilde{\mathbf{x}})}{d\tilde{\mathbf{x}}} = \mathbf{A}^\top + \mathbf{A} + 2\tilde{\mathbf{x}} + \mathbf{A}\tilde{\mathbf{x}} + \mathbf{A}^\top \tilde{\mathbf{x}}$

$$\nabla_{\tilde{\mathbf{x}}}(\tilde{\mathbf{y}} \cdot \tilde{\mathbf{z}}) = (\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{y}}) \tilde{\mathbf{z}} + (\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{z}}) \tilde{\mathbf{y}} \qquad \nabla_{\tilde{\mathbf{x}}} f(\tilde{\mathbf{y}}) = (\nabla_{\tilde{\mathbf{x}}} \tilde{\mathbf{y}})(\nabla_{\tilde{\mathbf{y}}} f(\tilde{\mathbf{y}}))$$
$$\nabla_w w^T A w = (A + A^T) w \qquad \mathbf{H}_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

Perceptron (Lecture 2)

$f(\tilde{\mathbf{x}}) = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}} + \alpha = \sum_{i=1}^d w_i x_i + \alpha$,
Goal: find w s.t all constraints $y_i X_i \cdot w \geq 0$. Define a risk function and optimize it, where the loss is defined as $L(z, y_i) = -y_i z$ if $y_i z < 0$, else 0. Therefore risk $R(w) = \sum_{i \in V} -y_i X_i \cdot w$

Decision boundary, a **hyperplane** in \mathbb{R}^d :
 $H = \{\tilde{\mathbf{x}} \in \mathbb{R}^d : f(\tilde{\mathbf{x}}) = 0\} = \{\tilde{\mathbf{x}} \in \mathbb{R}^d : \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}} + \alpha = 0\}$

$\tilde{\mathbf{w}}$ is the **normal** of the hyperplane,
 α is the **offset** of the hyperplane from origin,
 $\frac{f(\tilde{\mathbf{x}})}{\|\tilde{\mathbf{w}}\|}$ is the **signed distance** from the $\tilde{\mathbf{x}}$ to hyperplane \mathcal{H} .

Perceptron algorithm,
Input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{\pm 1\}$
while some $y_i \neq \text{sign}(\tilde{\mathbf{w}} \cdot \mathbf{x}_i)$
 pick some misclassified (\mathbf{x}_i, y_i)
 $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} + y_i \mathbf{x}_i$

Given a **linearly separable data**, perceptron algorithm will take no more than $\frac{R^2}{\gamma^2}$ updates to **converge**, where $R = \max_i \|\mathbf{x}_i\|$ is the radius of the data,
 $\gamma = \min_i \frac{y_i(\tilde{\mathbf{w}} \cdot \mathbf{x}_i)}{\|\tilde{\mathbf{w}}\|}$ is the margin.
Also, $\frac{\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}}{\|\tilde{\mathbf{w}}\|}$ is the signed distance from H to $\tilde{\mathbf{x}}$ in the direction $\tilde{\mathbf{w}}$.

Gradient descent view of perceptron, minimize margin cost function
 $J(\tilde{\mathbf{w}}) = \sum_i (-y_i(\tilde{\mathbf{w}} \cdot \mathbf{x}_i))_+$ with $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} - \eta \nabla J(\tilde{\mathbf{w}})$

Support Vector Machine (Lecture 3, 4)

Hard margin SVM,
This method makes the margin as wide as possible. The signed distance from the hyperplane to X_i is $\frac{f(\mathbf{x}_i)}{\|\tilde{\mathbf{w}}\|}$. Hence the margin is $\min_i \frac{1}{\|\tilde{\mathbf{w}}\|} |w \cdot X_i + \alpha| \geq \frac{1}{\|\tilde{\mathbf{w}}\|} \implies \min_{\tilde{\mathbf{w}}} \|\tilde{\mathbf{w}}\|^2$, such that $y_i \tilde{\mathbf{w}} \cdot \mathbf{x}_i \geq 1 (i = 1, \dots, n)$
Soft margin SVM,
 $\min_{\tilde{\mathbf{w}}} \|\tilde{\mathbf{w}}\|^2 + C \sum_{i=1}^n \xi_i$

Regularization and SVMs: Simulated data with many features $\phi(\tilde{\mathbf{x}})$; C controls trade-off between margin $1/\|\tilde{\mathbf{w}}\|$ 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.

Decision Theory (Lecture 6)

Bayes Theorem: $\underbrace{\mathbb{P}(Y = C|X)}_{\text{Poster. Prob}} = \frac{\overbrace{\mathbb{P}(X|Y=C)\mathbb{P}(Y=C)}^{\text{Prior Prob}}}{\mathbb{P}(X)}$
Assume **(X, Y)** are chosen i.i.d according to some probability distribution on $\mathcal{X} \times \mathcal{Y}$. **Risk** is misclassification probability:

$$R(r) = \mathbb{E}[L(r(\mathbf{X}), \mathbf{Y})] = \mathbb{P}(r(\mathbf{X}) \neq \mathbf{Y})$$
$$= \sum_{\tilde{\mathbf{x}}} [L(r(\tilde{\mathbf{x}}), 1) \mathbb{P}(Y = 1|x) + L(r(x), -1) \mathbb{P}(Y = -1|X = \tilde{\mathbf{x}})] \times \mathbb{P}(\tilde{\mathbf{x}})$$
$$= \mathbb{P}(Y = 1) \sum_x L(r(\tilde{\mathbf{x}}), 1) \mathbb{P}(\tilde{\mathbf{x}}|Y = 1) + \mathbb{P}(Y = -1) \sum_x L(r(\tilde{\mathbf{x}}), -1) \mathbb{P}(\tilde{\mathbf{x}}|Y = -1)$$

Bayes Decision Rule is
$$r^*(x) = \begin{cases} 1, & \text{if } L(-1, 1) \mathbb{P}(\mathbf{Y} = 1|x) > L(1, -1) \mathbb{P}(\mathbf{Y} = -1|x) \\ -1, & \text{otherwise.} \end{cases}$$
,
and the optimal risk (Bayes risk) $R^* = \inf_r R(r) = R(r^*)$

Generative and Discriminative Models (Lecture 6)

Discriminative models: $\mathbb{P}(\mathbf{X}, \mathbf{Y}) = \mathbb{P}(\mathbf{X}) \mathbb{P}(\mathbf{Y}|\mathbf{X})$.
Estimate $\mathbb{P}(\mathbf{Y}|\mathbf{X})$, then pretend our estimate $\hat{\mathbb{P}}(\mathbf{Y}|\mathbf{X})$ is the actual $\mathbb{P}(\mathbf{Y}|\mathbf{X})$ and plug in bayes rule expression.

Generative model: $\mathbb{P}(\mathbf{X}, \mathbf{Y}) = \mathbb{P}(\mathbf{Y}) \mathbb{P}(\mathbf{X}|\mathbf{Y})$.
Estimate $\mathbb{P}(\mathbf{Y})$ and $\mathbb{P}(\mathbf{X}|\mathbf{Y})$, then use bayes theorem to calculate $\mathbb{P}(\mathbf{Y}|\mathbf{X})$ and use discriminative model.

Gaussian class conditional densities $\mathbb{P}(\mathbf{X}|\mathbf{Y} = +1), \mathbb{P}(\mathbf{X}|\mathbf{Y} = -1)$ (with the same variance), the posterior probability is **logistic**:
 $\mathbb{P}(Y = +1|\tilde{\mathbf{x}}) = \frac{1}{1 + \exp(-\tilde{\mathbf{x}} \cdot \tilde{\mathbf{w}} - \beta_0)}$,
 $\tilde{\mathbf{w}} = \Sigma^{-1}(\mu_1 - \mu_0), \beta_0 = \frac{\mu_0' \Sigma^{-1} \mu_0 - \mu_1' \Sigma^{-1} \mu_1}{2} + \log \frac{\mathbb{P}(Y=1)}{\mathbb{P}(Y=0)}$

Multivariate Normal Distribution (Lecture 7)

$\tilde{\mathbf{x}} \in \mathbb{R}^d : p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{(-\frac{1}{2}(\tilde{\mathbf{x}} - \mu)' \Sigma^{-1}(\tilde{\mathbf{x}} - \mu))}$

QDA: Class-conditional densities $X_C \sim \mathcal{N}(\tilde{\mu}_C, \Sigma_C)$. Optimal decision rule $r^*(x)$ for 0–1 loss: Choose class **C** that maxes $\mathbb{P}(Y = C|X) \propto f_C(x)\pi_C$. Parameters estimated via MLE:
LDA: Assumes equal covariance matrices across classes ($\Sigma_C = \Sigma$), simplifying to linear decision surfaces.

$\Sigma = \mathbb{E}(\mathbf{X} - \mu)(\mathbf{X} - \mu)'$
Symmetric: $\Sigma_{i,j} = \Sigma_{j,i}$
Non-negative diagonal entries: $\Sigma_i, i \geq 0$
Positive semidefinite: $\forall \mathbf{v} \in \mathbb{R}^d, \mathbf{v}' \Sigma \mathbf{v} \geq 0$

Given a d -dimensaional Gaussian $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$, matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ and vector $\mathbf{b} \in \mathbb{R}^m$, define $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$. Then $\mathbf{Y} \sim \mathcal{N}(\mathbf{A}\mu + \mathbf{b}, \mathbf{A}\Sigma\mathbf{A}^\top)$

Given a d -dimensaional Gaussian $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$, with Σ positive definite,
 $\mathbf{Y} = \Sigma^{-\frac{1}{2}}(\mathbf{X} - \mu) \sim \mathcal{N}(\tilde{\mathbf{0}}, \mathbf{I})$

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

Prior: $\hat{\pi}_C = \mathbb{P}(Y = C) = \frac{N_C}{n}$
Mean: $\hat{\mu}_C = \mathbb{E}[\mathbf{X}|Y = C] = \frac{1}{N_C} \sum_{i: Y_i=C} X_i$

Covariance: $\hat{\Sigma}_C = \frac{1}{N_C} \sum_{i: Y_i=C} (X_i - \hat{\mu}_C)(X_i - \hat{\mu}_C)^\top$

Pooled Cov: $\hat{\Sigma} = \frac{1}{n} \sum_k \sum_{i: Y_i=C_k} (X_i - \hat{\mu}_{C_k})(X_i - \hat{\mu}_{C_k})^\top$

Discriminant Analysis (Lecture 7)

Discriminant Functions for LDA and QDA for class C , denoted $Q_C(\tilde{\mathbf{x}})$ is:
$$Q_C(\tilde{\mathbf{x}}) = \ln \left(\frac{f_{\tilde{\mathbf{X}}|Y=C}(\tilde{\mathbf{x}})\pi_C}{(2\pi)^{\frac{d}{2}}} \right) = -\frac{1}{2}(\tilde{\mathbf{x}} - \mu_C)^T \Sigma_C^{-1}(\tilde{\mathbf{x}} - \mu_C) - \frac{1}{2} \ln |\Sigma_C| + \ln \frac{\pi_C}{(2\pi)^{\frac{d}{2}}}$$

Here we have, the class-conditional density $f_{\tilde{\mathbf{X}}|Y=C}(\tilde{\mathbf{x}})$, prior probability π_C , the mean vector of class C , and covariance matrix Σ_C of class C .

Linear Decision Function between Classes C and D: Compares classes by:
$$Q_C(\tilde{\mathbf{x}}) - Q_D(\tilde{\mathbf{x}}) = (\mu_C - \mu_D)^T \Sigma^{-1} \tilde{\mathbf{x}} - \frac{1}{2}(\mu_C^T \Sigma^{-1} \mu_C - \mu_D^T \Sigma^{-1} \mu_D) + \ln \frac{\pi_C}{\pi_D}.$$

This expression includes a linear term $(\mu_C - \mu_D)^T \Sigma^{-1} \tilde{\mathbf{x}}$, which describes how the decision boundary is formed linearly in the feature space, and a constant term that adjusts the boundary based on the means of the classes and their prior probabilities.

Multi-class LDA: The optimal class for a given feature vector $\tilde{\mathbf{x}}$ is chosen by
Choose class C such that $Q_C(\tilde{\mathbf{x}})$ is maximized.

Effectively determining the most likely class based on the feature measurements and statistical properties of each class.

Linear Regression (Lecture 10)

Objective Function: The goal in linear regression is to minimize the sum of squared residuals (RSS), expressed as:

$$\text{RSS}(\tilde{\mathbf{w}}) = \sum_{i=1}^n (\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{w}} + \alpha - y_i)^2 = \|\mathbf{X}\tilde{\mathbf{w}} - \tilde{\mathbf{y}}\|^2.$$

Where, $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$ is the **design matrix**, and $\tilde{\mathbf{y}}$ is the **response vector**.
Minimize via Calculus: $\nabla \text{RSS}(\tilde{\mathbf{w}}) = 2\mathbf{X}^\top \mathbf{X}\tilde{\mathbf{w}} - 2\mathbf{X}^\top \tilde{\mathbf{y}}$.
Setting this gradient to zero leads to the **normal equations**:

$$\mathbf{X}^\top \mathbf{X}\tilde{\mathbf{w}} = \mathbf{X}^\top \tilde{\mathbf{y}} \implies \tilde{\mathbf{w}}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \tilde{\mathbf{y}},$$

$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is known as the *Moore-Penrose pseudo-inverse* X^\dagger of \mathbf{X} .
Projection Theorem also leads to the normal equations, by asserting that the solution $\tilde{\mathbf{w}}^*$ projects the residuals orthogonally onto the column space of \mathbf{X} :
$$\mathbf{X}^\top (\tilde{\mathbf{y}} - \mathbf{X}\tilde{\mathbf{w}}) = 0 \implies \mathbf{X}^\top \mathbf{X}\tilde{\mathbf{w}} = \mathbf{X}^\top \tilde{\mathbf{y}},$$

Logistic Regression (Lecture 10, 11)

Fits “probabilities” to data and usually used for classification.
 $\mathbb{P}(Y = 1|\tilde{\mathbf{x}}_i) = \frac{1}{1 + \exp(-\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_i)} = \sigma(\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_i)$
Logistic Loss: $L(z, y) = -y \ln z - (1 - y) \ln(1 - z)$

Cost Function: $J(\vec{w}) = -\sum_{i=1}^n L(\sigma(\vec{w}^\top \vec{x}_i), y_i)$

Find $\vec{w}^* = \text{argmin } J(\vec{w})$, Cost function is convex, and solved by gradient descent. Converges to max margin classifier.

Let $\sigma_i = \sigma(\vec{w}^\top \vec{x}_i)$ and $\vec{\sigma} = [\sigma_1 \quad \dots \quad \sigma_n]^\top$

$$\begin{aligned}\nabla_{\vec{w}} \sigma_i &= \sigma_i(1 - \sigma_i) \vec{x}_i \\ \nabla_{\vec{w}} J(\vec{w}) &= \mathbf{X}^\top (\vec{\sigma} - \vec{y}) \\ \nabla_{\vec{w}}^2 J(\vec{w}) &= \mathbf{X}^\top \text{diag}(\vec{\sigma}(1 - \vec{\sigma})) \mathbf{X}\end{aligned}$$

Gradient descent:

$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \nabla_{\vec{w}} J(\vec{w}^{(t)}) : O(nd)$ per step

Stochastic gradient descent:

$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \nabla_{L_i} J(\vec{w}^{(t)}) : O(d)$ per step

Newton's method:

$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \left(\nabla^2 J(\vec{w}^{(t)}) \right)^{-1} \nabla J(\vec{w}^{(t)})$

LDA vs. Logistic Regression

Advantages of LDA:

- For well-separated classes, LDA is stable; log. reg. is surprisingly unstable.
- > 2 classes? Easy & elegant; log. reg. needs modifying. (softmax regr.)
- slightly more accurate when classes nearly normal, especially if n is small.

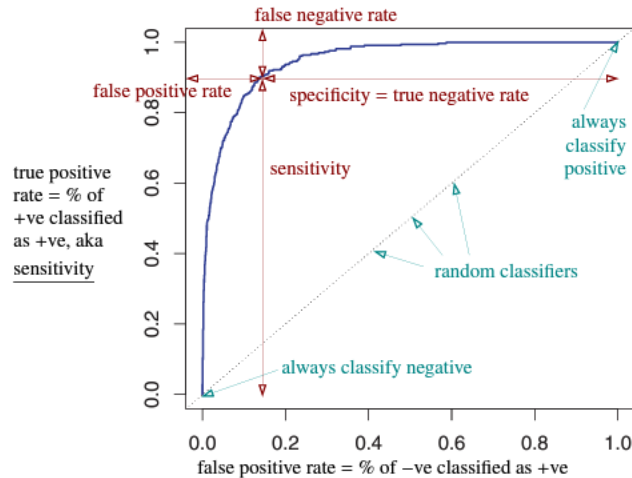
Advantages of Log. Regression:

More emphasis on desc. boundary; always separates linearly separable pts.

More robust on some non-Gaussian distributions (e.g., dists w/ large skews.)

Naturally fits labels between 0 and 1.

ROC CURVES



Polynomial Regression (Lecture 11)

Replace each x_i with feature vector $\Phi(x_i)$ with all terms of degree p

polynomial. e.g., $\Phi(x_i) = [x_{i1}^2 \quad x_{i1}x_{i2} \quad x_{i2}^2 \quad x_{i1} \quad x_{i2} \quad 1]^\top$

Log. Reg. + quadratic features = same form of posteriors as QDA.

Statistical Justification For Regression (Lecture 12)

Typical model of reality: $\forall X_i, y_i = g(X_i) + \epsilon_i : \epsilon_i \sim D'$ where D' has mean 0.

Ideal approach: choose $h(X_i) = \mathbb{E}_Y[Y|X = X_i] = g(X_i) + \mathbb{E}[\epsilon_i] = g(X_i)$

Least Squares Cost Function From MLE

Suppose $\epsilon_i \sim \mathcal{N}(0, \sigma^2) \Rightarrow y_i \sim \mathcal{N}(g(\vec{x}_i), \sigma^2)$. We're going to try to estimate $g(\vec{x})$, which is defined by some weights \vec{w} . The probability of y_i , given it's parameters for its distribution, is a pdf $f(y_i)$ & the log likelihood is:

$$\ell(g; \mathbf{X}, \vec{y}) = \ln \left(\prod_{i=1}^n f(y_i) \right) = \sum_{i=1}^n \ln(f(y_i)) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - g(\vec{x}_i))^2 - C$$

Maximizing this is equivalent to minimizing $\sum_{i=1}^n (y_i - g(\vec{x}_i))^2 = \text{RSS}$

Logistic Loss from MLE

Consider $P(\vec{x}_i \in \text{class C}) = y_i$ as the actual probability, with $h(\vec{x}_i)$ as the predicted probability. For β hypothetical repetitions, where $y_i \beta$ samples belong to class C and $(1 - y_i)\beta$ do not. The likelihood of h given the data is:

$$\mathcal{L}(h; \mathbf{X}, \vec{y}) = \prod_{i=1}^n h(\vec{x}_i)^{y_i \beta} (1 - h(\vec{x}_i))^{(1 - y_i) \beta}$$

$$\ell(h; \mathbf{X}, \vec{y}) = \beta \sum_{i=1}^n [y_i \ln(h(\vec{x}_i)) + (1 - y_i) \ln(1 - h(\vec{x}_i))]$$

Maximizing the log likelihood $\ell(h)$ is equivalent to minimizing logistic losses.

The Bias-Variance Decomposition

There are 2 sources of error in a hypothesis h :

bias: error due to inability of hypothesis h to fit g perfectly.

variance: error due to fitting random noise in data, e.g., we fit linear g with a linear h , yet $h \neq g$.

Model: $X_i \sim D, \epsilon_i \sim D'$ and $y_i = g(X_i) + \epsilon_i$. Fit hypothesis h to \mathbf{X}, \vec{y} . Now h is a r.v., because of the random noise ϵ . Let z be an arbitrary pt and $\gamma = g(z) + \epsilon$. Now the risk fn when loss is the squared error is:

$R(h) = \mathbb{E}[L(h(z), \gamma)] \leftarrow \text{exp over possible training set } X, y \text{ and values of } \gamma$

$$R(h) = \mathbb{E}[(h(z) - \gamma)^2]$$

$$R(h) = \mathbb{E}[h(z)^2] + \mathbb{E}[\gamma^2] - 2\mathbb{E}[\gamma h(z)]$$

$$R(h) = \text{Var}(h(z)) + \mathbb{E}[h(z)]^2 + \text{Var}(\gamma) + \mathbb{E}[\gamma]^2 - 2\mathbb{E}[\gamma] \mathbb{E}[h(z)]$$

$$R(h) = (\mathbb{E}[h(z)] - \mathbb{E}[\gamma])^2 + \text{Var}(h(z)) + \text{Var}(\gamma)$$

$$R(h) = \underbrace{(\mathbb{E}[h(z)] - \mathbb{E}[\gamma])^2}_{\text{bias}^2 \text{ of method}} + \underbrace{\text{Var}(h(z))}_{\text{variance of method}} + \underbrace{\text{Var}(\gamma)}_{\text{irreducible error}}$$

- **Underfitting:** Too much bias
- **Overfitting:** Too much variance
- Training Error reflects bias but not variance
- Test error reflects both.

Linear Regression Regularization (Lecture 13)

Trading off bias/variance: some increase in bias can give big decrease in var.

Ridge regression is like L_2 regularization:

$\hat{\vec{w}} = \arg \min_{\vec{w}} \|\vec{y} - \mathbf{X}\vec{w}\|^2 + \lambda \|\vec{w}\|^2$, for some cost function J on classifier h

Minimized w/ Calculus

$$\hat{\vec{w}}^{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \vec{y}$$

Lasso regression is like L_1 regularization:

Minimized w/ Calculus (it's also a quadratic program)

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

Misc

Centering X: This involves subtracting μ^T from each row of \mathbf{X} . Symbolically, \mathbf{X} transforms into $\tilde{\mathbf{X}}$.

Decorrelating X: This process applies a rotation $\mathbf{Z} = \tilde{\mathbf{X}}\mathbf{V}$, where $\text{Var}(\mathbf{R}) = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$. This step rotates the sample points to the eigenvector coordinate system.

Sphering: $\tilde{\mathbf{X}}$: applying transform $\mathbf{W} = \tilde{\mathbf{X}}\text{Var}(\mathbf{R})^{-\frac{1}{2}}$
whitening X: centering + sphering, $\mathbf{X} \rightarrow \mathbf{W}$

Decision Trees (Lecture 14, 15)

Nonlinear method for classification or regression. Cuts x-space into rectangular cells. Works well with both categorical and quantitative features. They're fast, interpretable, robust to irrelevant features, but not best at prediction \rightarrow high variance often mitigated by ensemble learning.

For classification the learning algorithm is a greedy, top-down learning heuristic. Let S be subset of sample point indices:

function GROWTREE(S)

if all $y_i = C$ for all $i \in S$ and some class C **then**

return new leaf(C)

else

Choose best splitting feature k and splitting value β

$S_L = \{i \in S : \mathbf{X}_{ik} < \beta\}$

$S_R = \{i \in S : \mathbf{X}_{ik} \geq \beta\}$

return new node(k, β , GROWTREE(S_L), GROWTREE(S_R))

Optimal Split Selection: Evaluate all possible (k, β) splits and choose the one that minimizes the weighted impurity cost:

$$J(S, k, \beta) = \frac{|S_L|J(S_L) + |S_R|J(S_R)}{|S_L| + |S_R|}$$

where entropy $J(S) = -\sum_c p_c \log_2(p_c)$, with $p_c = \frac{|\{i \in S : y_i = c\}|}{|S|}$

Information Gain, defined as $IG(S, k, \beta) = J(S) - J(S, k, \beta)$

IG maximization is equivalent to entropy minimization post-split.

Pruning: Grow tree too large; greedily remove each split whose removal improves validation performance on test data that lies on these splits.

For regression: Leafs stores mean response for training pts in that region.

Different cost function: $J(S) = \text{RSS} = \frac{1}{|S|} \sum_{i \in S} (y_i - \mu_S)^2 = \text{Var}(\{y_i : i \in S\})$
Ensemble Methods Decision trees often suffer from high variance, which can lead to overfitting. Ensemble methods, by combining multiple trees, aim to reduce this variance and improve predictive performance.

Bagging: Given n -point training samples, generate random subsamples of size n' by sampling with replacement. If $n' = n \sim 63.2\%$ are chosen. Idea: Points chosen j times have greater weight in cost when misclassified.

Random Forest: Enhance Bagging by reducing correlation between different trees with feature randomness at each split. This process ensures that the ensemble does not rely too heavily on any single feature and makes trees within the forest less identical.

Kernels (Lecture 16)

In many learning algorithms, the weights can be written as a linear combo of training points, that is, $\tilde{\mathbf{w}} = \mathbf{X}^\top \tilde{\boldsymbol{\alpha}} = \sum_{i=1}^n \alpha_i \tilde{\mathbf{x}}_i$, for some $\tilde{\boldsymbol{\alpha}} \in \mathbb{R}^n$.

We can substitute this identity into alg, and optimize n dual weights α instead of the $\tilde{\mathbf{w}}$ of the primal problem. **Kernel Trick:** Enables learning algorithms to operate in high-dimensional space without explicitly computing the transformations.

Kernel Ride Regression

Transform the training set X using a kernel function $k(x, z)$, creating a kernel matrix K .

Solve $(K + \lambda I)a = y$ for the dual coefficients a , where λ is the regularization parameter.

Predictions for a new point z are made by $h(z) = \sum_{i=1}^n \alpha_i k(x_i, z)$.

Center \mathbf{X} : $\tilde{\mathbf{x}}_i \leftarrow \tilde{\mathbf{x}}_i - \mu_x$, except fictitious dimension, $\tilde{\mathbf{x}}_{i,d+1} = 1$

Center $\tilde{\mathbf{y}}$: $\tilde{y} \leftarrow \tilde{y} - \mu_y$.

In ridge regression, recall the normal equations:

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}') \tilde{\mathbf{w}} = \mathbf{X}^\top \tilde{\mathbf{y}}$$

Where \mathbf{I}' is the identity matrix with 0 at the position of the fictitious dimension so that we don't penalize the bias. In order for the condition that the weights be a linear combination of training points to be satisfied, we need to penalize the bias.

Now the objective function is: (subsiting $\tilde{\mathbf{w}} = \mathbf{X}^\top \tilde{\boldsymbol{\alpha}}$)

$$\left\| \mathbf{X}^\top (\mathbf{X} \tilde{\boldsymbol{\alpha}}) - \tilde{\mathbf{y}} \right\|^2 + \lambda \left\| \mathbf{X}^\top \tilde{\boldsymbol{\alpha}} \right\|^2$$

Training: Solve $(\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}) \tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{y}}$, for $\tilde{\boldsymbol{\alpha}}$.

Testing: Regression fn is $h(\tilde{\mathbf{z}}) = \tilde{\mathbf{w}}^\top \tilde{\mathbf{z}} = \tilde{\boldsymbol{\alpha}}^\top \mathbf{X} \tilde{\mathbf{z}} = \sum_{i=1}^n \alpha_i (\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{z}})$

Let $k(\tilde{\mathbf{x}}, \tilde{\mathbf{z}}) = \tilde{\mathbf{x}}^\top \tilde{\mathbf{z}}$ be the kernel function.

Let $\mathbf{K} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{n \times n}$ be the kernel matrix with $\mathbf{K}_{ij} = k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)$.

Kernel Ridge Regr. Alg:

Initialize the kernel matrix \mathbf{K} $\triangleright \mathcal{O}(n^2 d)$ time
Solve $(\mathbf{K} + \lambda \mathbf{I}) \tilde{\boldsymbol{\alpha}} = \tilde{\mathbf{y}}$ $\triangleright \mathcal{O}(n^3)$ time
for each test point z , compute **do**
 $h(z) = \sum_{i=1}^n \alpha_i k(\tilde{\mathbf{x}}_i, z)$ $\triangleright \mathcal{O}(nd)$ time per test pt

The Kernel Trick (aka Kernelization)

Polynomial kernel of degree p is $k(\tilde{\mathbf{x}}, \tilde{\mathbf{z}}) = (\tilde{\mathbf{x}}^\top \tilde{\mathbf{z}} + 1)^p$

Theorem: $(\tilde{\mathbf{x}}^\top \tilde{\mathbf{z}} + 1)^p = \Phi^T(\tilde{\mathbf{x}}) \Phi(\tilde{\mathbf{z}})$ for some feature map Φ

Note now we can compute $\Phi^T(\tilde{\mathbf{x}}) \Phi(\tilde{\mathbf{z}})$ in $\mathcal{O}(d)$ time instead of $\mathcal{O}(d^p)$

Gaussian (RBF) Kernel: $k(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$, measures the similarity, or closeness, of x and z with a parameter σ controlling the width of the Gaussian.

Choosing σ in Gaussian Kernels:

Trade-off between bias and variance: larger σ leads to smoother models with higher bias and lower variance.

Selected based on validation performance, optimizing for generalization.

Neural Networks (Lecture 17)

Input Layer: $x_1, \dots, x_d : x_{d+1} = 1$

Hidden Units: $h_1, \dots, h_m : h_{m+1} = 1$

Output Layer: z_1, \dots, z_k

Layer 1 weights $V \in \mathbb{R}^{m \times (d+1)}$ (row i , V_i^\top contains all the weights from previous layer, coming into hidden unit i)

Layer 2 weights $W \in \mathbb{R}^{k \times (m+1)}$

$\tilde{\mathbf{z}} = \tilde{\mathbf{s}}(\mathbf{W}\tilde{\mathbf{h}}) = \tilde{\mathbf{s}}(\mathbf{W}\tilde{\mathbf{s}}_1(\mathbf{V}\tilde{\mathbf{x}}))$ where $\tilde{\mathbf{s}}_1$ is the activation function applied element wise to vector, with last component set to 1.

Activation Functions:

$$s(\gamma) = \sigma(\gamma) = \frac{1}{1 + e^{-\gamma}} \implies \sigma'(\gamma) = \sigma(\gamma)(1 - \sigma(\gamma))$$

$$s(\gamma) = \sigma_{\text{ReLU}}(\gamma) = \max(0, \gamma) \implies \sigma'_{\text{ReLU}}(\gamma) = \begin{cases} 1 & \gamma > 0 \\ 0 & \gamma \leq 0 \end{cases}$$

- Derive the gradients of the loss $L \in \mathbb{R}$ with respect to weight matrix $W \in \mathbb{R}^{n^{(l)} \times n^{(l+1)}}$ and bias row vector $b \in \mathbb{R}^{1 \times n^{(l+1)}}$ in the fully-connected layer, $\partial L / \partial W$ and $\partial L / \partial b$. You will also need to take the gradient of the loss with respect to the input of the layer $\partial L / \partial X$, which will be passed to lower layers, where $X \in \mathbb{R}^{m \times n^{(l)}}$ is the batched input. **Again, you must arrive at a solution that uses batched X and Z . Please express your solution in terms of $\partial L / \partial Z$, which you have already obtained in 4.1, where $Z = XW + 1b$ and $1 \in \mathbb{R}^{m \times 1}$ is a column of ones.** Note that $1b$ is a matrix (it's an outer product) whose each row is the row vector b so we are adding the same bias vector to each sample during the forward pass: this is the mathematical equivalent of numpy broadcasting.

Solution: Let L be a function operating on Z , where $Z = XW + 1b$. Then:

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial W} = X^\top \frac{\partial L}{\partial Z} \quad (1)$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial b} = 1^\top \frac{\partial L}{\partial Z} \quad (2)$$

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Z} \frac{\partial Z}{\partial X} = \frac{\partial L}{\partial Z} W^\top \quad (3)$$

- First, derive the gradient of the downstream loss with respect to the input of the ReLU activation function, Z . **You must arrive at a solution in terms of $\partial L / \partial Y$, the gradient of the loss w.r.t. the output of ReLU $Y = \sigma_{\text{ReLU}}(Z)$, and the batched input Z , i.e., where $Z \in \mathbb{R}^{m \times n}$.**

Hint: you are allowed to use operations like elementwise multiplication and/or division!

Solution: We know that $\sigma'_{\text{ReLU}}(\gamma) = \begin{cases} 0 & \gamma < 0, \\ 1 & \text{otherwise.} \end{cases}$

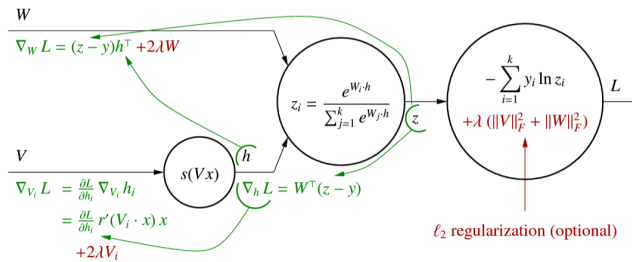
L is some function operating on Y , the output of ReLU, that is:

$$L = f(Y) = f(\sigma_{\text{ReLU}}(Z)).$$

Thus, by chain rule, the gradient of L with respect to Z is.

$$\begin{aligned} \frac{\partial L}{\partial Z} &= \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial Z} \\ &= \frac{\partial L}{\partial Y} \circ \sigma'_{\text{ReLU}}(Z) \end{aligned}$$

Where \circ is the element-wise multiplication.



PCA: Principal Component Analysis (Lecture 20)

Goal: Identify k principal directions to maximize variance capture after projecting data onto these components, enhancing computational efficiency and reducing overfitting by eliminating irrelevant features.

Assumptions: Assume \mathbf{X} is centered, i.e., $\sum_{i=1}^n \mathbf{x}_i = 0$.

Orthogonal Projection of a point \mathbf{x} onto a vector \mathbf{w} is given by:

$$\tilde{\mathbf{x}} = \frac{\mathbf{x} \cdot \mathbf{w}}{\|\mathbf{w}\|^2} \mathbf{w}$$

Aim: Find the best \mathbf{w} that captures the maximum variance after projecting.

Principal Coordinates and Components: Given orthonormal directions $\mathbf{v}_1, \dots, \mathbf{v}_k$, the projection of \mathbf{x} onto these vectors is:

$$\tilde{\mathbf{x}} = \sum_{i=1}^k (\mathbf{x} \cdot \mathbf{v}_i) \mathbf{v}_i$$

where $\mathbf{x} \cdot \mathbf{v}_i$ are the principal coordinates. The principal components are unit eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_d$ of the covariance matrix $\hat{\Sigma} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$, associated with its eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$.

function PCA ALGORITHM($\mathbf{X}, \mathbf{X}_{\text{test}}$)

Center \mathbf{X}

Compute unit eigenvectors and eigenvalues of $\mathbf{X}^\top \mathbf{X}$

Choose k (optionally based on the size of eigenvalues)

Select the top k eigenvectors $\mathbf{v}_{d-k+1}, \dots, \mathbf{v}_d$ for the best k -dimensional subspace

Compute the k principal coordinates $\mathbf{x} \cdot \mathbf{v}_i$ for each training/test point

Derivation 2: Variance Maximization: Find \mathbf{w} that maximizes

$$\text{Var}(\{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_n\}) = \frac{1}{n} \frac{\|\mathbf{X}\mathbf{w}\|^2}{\|\mathbf{w}\|^2} = \frac{1}{n} \frac{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} = \frac{1}{n} \text{Rayleigh Quotient}(\mathbf{X}^\top \mathbf{X})$$

Theorem: $\tilde{\mathbf{v}}_d$ maximizes this variances.

Derivation 3: Projection Distance Minimization:

Find direction \mathbf{w} that minimizes the mean of squared projection distance.

$$\begin{aligned} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 &= \sum_{i=1}^n \left\| \mathbf{x}_i - \frac{\mathbf{x}_i \cdot \mathbf{w}}{\|\mathbf{w}\|} \mathbf{w} \right\|^2 \\ &= \sum_{i=1}^n \left(\|\mathbf{x}_i\|^2 - \left(\mathbf{x}_i \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^2 \right) = C - n(\text{variation from derivation 2}) \end{aligned}$$

Minimizing mean squared project distance \equiv Maximizing variance.

SVD, Clustering (Lecture 21)

SVD: Problem with PCA is that computing $\mathbf{X}^\top \mathbf{X}$ is expensive, taking $\mathcal{O}(nd^2)$.

$$[\mathbf{X}]_{n \times d} = [\mathbf{U}]_{n \times d} [\mathbf{D}]_{d \times d} [\mathbf{V}]_{d \times d}^\top = \sum_{i=1}^d \delta_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^\top \quad \text{when } n \geq d$$

- The number of non-negative singular values corresponds to the rank of \mathbf{X} .
- $\tilde{\mathbf{v}}_i$ is an eigenvector of $\mathbf{X}^\top \mathbf{X}$ with eigenvalue δ_i^2 .
- We can find the k greatest singular values & corresp vectors in $\mathcal{O}(ndk)$
- Row i of $\mathbf{U}\mathbf{D}$ gives the principle coordinates of sample \mathbf{X}_i

Clustering: Partition data into K clusters so pts in a cluster are more similar than across clusters.

K-Means Clustering goal is to partition pts into k disjoint clusters to minimize the sum of squared distances from each point to its cluster center.

$$\text{Find } \tilde{\mathbf{y}} \text{ that minimizes } \sum_{i=1}^k \sum_{y_j=i} \left\| \mathbf{X}_j - \tilde{\boldsymbol{\mu}}_i \right\|^2$$

Where, $\tilde{\boldsymbol{\mu}}_i = \frac{1}{n_i} \sum_{y_j=i} \mathbf{X}_j$ is the mean of cluster i , given n_i pts in cluster i .

Problem is NP-hard, $\mathcal{O}(nk^n)$. Use heuristic, alternating between 2 steps.

- Fixing labels y_j and updating cluster means μ_i to min the obj. fn.
- Fixing $\tilde{\boldsymbol{\mu}}_i$ and updating labels y_j to minimize the objective function.

Repeats until no assignments change. Both steps decrease obj. fn, *unless* they change nothing. Hence alg. must terminate. We can show through calculus:

- The optimal $\tilde{\boldsymbol{\mu}}_i$ is the mean of its cluster

2. The optimal \tilde{y} assigns each pt \mathbf{x}_j to nearest center $\tilde{\mu}_i$

Initialization Strategies

- Forgy Method:** Choose k random sample pts to be initial $\tilde{\mu}_i$'s; go to step 2.
- Random Partition:** randomly assign each sample pt to a cluster; go to step 1.
- K-Means++:** Like Forgy, but biased distribution, choosing centers far apart.

Equivalent Object Function: called *The Within-Cluster Variation*: (no $\tilde{\mu}_i$)

$$\min_{\tilde{y}} \sum_{i=1}^k \frac{1}{n_i} \sum_{y_j=i} \sum_{y_m=i} \left\| \mathbf{x}_j - \mathbf{x}_m \right\|^2$$
. Minimum is twice the prev. obj. fn.

Hierarchical Clust: Builds tree of clusters; every subtree is cluster. Built using

- Agglomerative (Bottom-up):** Start with each point as a cluster, then merge the closest pairs until one cluster remains.
- Divisive (Top-down):** Start with all points in one cluster, then iteratively split the most diverse cluster.

We need a distance fn for clusters A, B . Here are 4 common ones :

- Complete Linkage:** Max distance between any points in two clusters.
- Single Linkage:** Min distance between any points in two clusters.
- Average:** Avg dist between all points in clusters.
$$d(A, B) = \frac{1}{|A||B|} \sum_{w \in A} \sum_{x \in B} d(w, x)$$
- Centroid Linkage:** dist between centroids of clusters. $d(A, B) = d(\mu_A, \mu_B)$

Greedy Agglomerative Alg: repeatedly *fuse* the 2 clusters that minimize $d(A, B)$. Naive takes $\mathcal{O}(n^3)$, done in $\mathcal{O}(n^2)$ for complete and single linkage
Dendrogram: Illustration of the cluster hierarchy. Vertical axis encodes all the linkage distances. Horizontal axis has no meaning. The complete linkage tends to give balanced trees.

Learning Theory (Lecture 23)

Learning theory explores how machine learning algs generalize from training data to unseen data. It involves understanding the constraints needed for a learner to make predictions on new data effectively.

Range Space: Defined by (P, H) where P is the set of all train/test points and H is the hypothesis class or a set of classifiers. Each hypothesis $h \in H$ identifies which points are predicted to be in a particular class.

Risk and Empirical Risk: *Risk (Generalization Error) $R(h)$:* The prob. that a hypothesis h misclassifies a randomly drawn point from distribution \mathcal{D} .
Empirical Risk $\hat{R}(h)$: The proportion of training data misclassified by h , which better approximates $R(h)$ as the size of the training set increases.

Hoeffding's Inequality: Tells us probability of bad Risk estimate.

$$\mathbb{P}(|\hat{R}(h) - R(h)| > \epsilon) \leq 2e^{-2\epsilon^2 n},$$

where ϵ is the error margin and n is # of samples. This shows that larger n exponentially decays the chance of deviation between $\hat{R}(h)$ and $R(h)$.

Empirical Risk Minimization (ERM): Selects the hypothesis \hat{h} that minimizes $\hat{R}(\hat{h})$ among H . ERM is infeasible to pick optimal hypothesis.

Dichotomies and Sample Complexity: A dichotomy created by a hypothesis h over set X categorizes points into class C or its complement. The number of possible dichotomies, Π , influences the likelihood of misleadingly low empirical risks. If H induces Π dichotomies, then:

$$\mathbb{P}(\text{at least one dichotomy has } |\hat{R} - R| > \epsilon) \leq \delta = 2\Pi e^{-2\epsilon^2 n}$$

with confidence $1 - \delta$, $\forall h \in H$, $|\hat{R}(h) - R(h)| \leq \epsilon = \sqrt{\frac{1}{2n} \ln \frac{2\Pi}{\delta}}$. This formulation reveals that increasing n and reducing Π enhances risk estimates' reliability.

Optimal Risk Approximation: If $h^* \in H$ minimizes $R(h^*)$, then with high probability $1 - \delta$, the selected hypothesis \hat{h} that minimizes $\hat{R}(\hat{h})$ fulfills:

$$R(\hat{h}) \leq \hat{R}(\hat{h}) + \epsilon \leq \hat{R}(h^*) + 2\epsilon$$

Sample Complexity: The minimum # of samples required to ensure that a approximation accuracy with confidence $1 - \delta$ is:

$$n \geq \frac{1}{2\epsilon^2} \ln \frac{2\Pi}{\delta}$$

AdaBoost: Adaptive Boosting (Lecture 24)

An ensemble method for classification (or regression) that trains T learners on weighted samples points, gives bigger votes to more accurate learners, and increases weights of misclassified points (pay more attention to this point), with the goal of reducing bias. If every learner has $\text{err} < 0.5$, $M(\tilde{x})$ training $\text{err} \rightarrow 0$ as $T \rightarrow \infty$.

Input: $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\tilde{y} \in \mathbb{R}^n$ with $y_i = \pm 1$ (not natural fit if labels arent bin).
Idea:

- Training Classifiers:** Sequentially train T classifiers G_1, \dots, G_T
- Weight Adjustment:** Weight for training pt \tilde{x}_i in iteration t for G_t grows according to how many of G_1, \dots, G_{t-1} misclassified it.
- Meta Learner:** is a linear combination of all T classifiers. Output is not binary; apply the sign function for binary classification.

Risk Minimization: Risk of the meta learner is average exponential loss $L(\rho, \ell) = e^{-\rho \ell}$, and minimize it to adjust β_T and select the optimal classifier G_T for each round:

$$\begin{aligned} n * \text{Risk} &= \sum_{i=1}^n L(M(\tilde{x}_i), y_i) = \sum_{i=1}^n \prod_{t=1}^T \exp(-y_i \beta_t G_t(\tilde{x}_i)) \\ &= \sum_{i=1}^n w_i^{(T)} \exp(-y_i \beta_T G_T(\tilde{x}_i)) \\ &= e^{-\beta_T} \sum_{y_i = G_T(\tilde{x}_i)} w_i^{(T)} + e^{\beta_T} \sum_{y_i \neq G_T(\tilde{x}_i)} w_i^{(T)} \\ &= e^{-\beta_T} \sum_{i=1}^n w_i^{(T)} + \cancel{(e^{\beta_T} - e^{-\beta_T})} \sum_{y_i \neq G_T(\tilde{x}_i)} w_i^{(T)} \end{aligned}$$

where $w_i^{(T)} = \prod_{t=1}^{T-1} \exp(-y_i \beta_t G_t(\tilde{x}_i))$ is the cumulative weight adjusted for each sample up to classifier $T - 1$. Therefore, the optimal classifier G_T minimizes the sum of weights of misclassified points.

Recursive Weight Update: To focus training on misclassified points:

$$w_i^{(T+1)} = w_i^{(T)} \exp(-y_i \beta_T G_T(\tilde{x}_i)) = \begin{cases} w_i^{(T)} \exp(-\beta_T) & \text{if } y_i = G_T(\tilde{x}_i) \\ w_i^{(T)} \exp(\beta_T) & \text{if } y_i \neq G_T(\tilde{x}_i) \end{cases}$$

Optimization for β_T : Solve for β_T by setting the derivative of the risk with respect to β_T equal to zero, deriving that:

$$\beta_T = \frac{1}{2} \ln \left(\frac{1 - \text{err}_T}{\text{err}_T} \right),$$

where err_T is the weighted proportion of misclassified instances by G_T . Therefore a perfect learner ($\text{err}_T = 0$) has $\beta_T = \infty$

function ADABOOST(\mathbf{X}, \tilde{y})
 Initialize weights $w_i = \frac{1}{n}$ for all i
 for $t \leftarrow 1$ to T **do**
 Train classifier G_t on (\mathbf{X}, \tilde{y}) , using weights \tilde{w}
 Calculate error err_t
 Calculate β_t using err_t
 Update weights: $w_i \leftarrow w_i \exp(-y_i \beta_t G_t(x_i))$
 Return metalearner $h(\tilde{z}) = \text{sign}(\sum_{t=1}^T \beta_t)$

K-Nearest Neighbor (Lecture 25)

- Idea:** Given query point \tilde{q} , find the k training pts nearest to \tilde{q} .
- Theorem:** as $n \rightarrow \infty$, the 1-NN error rate is $< 2B - B^2$ where B = bayes risk.
- Theorem:** As $n \rightarrow \infty$, $k \rightarrow \infty$, and $\frac{k}{n} \rightarrow 0$ (n grows faster than k), k-NN error rate converges to B

Exhaustive k-NN: Algorithm:

Given query point q :

- Scan through all n training pts, computing (squared) distances to q .
- Maintain a max-heap with the k closest training points seen so far.
- Time to train classifier: $\mathcal{O}(0)$ (no time)
- Query Time: $\mathcal{O}(nd + n \log k)$, expected $\mathcal{O}(nd + k \log(n) \log(k))$ if random point order.

Voronoi Diagrams:

Let X be a point set. the **Voronoi cell** of $w \in X$ is
 $\text{Vor } w = \{p \in \mathbb{R}^d : \|p - w\| \leq \|p - v\| \ \forall v \in X\}$
The *Voronoi diagram* of X is the set of X 's Voronoi cells.

K-Nearest Neighbor (Lecture 25)

Principle: For a query point \tilde{q} , classify based on the majority label among its k nearest training points.

Theorem: Under certain conditions:

- For 1-NN, as $n \rightarrow \infty$, the error rate is less than twice the Bayes risk B .
- For k-NN, as $n \rightarrow \infty$, $k \rightarrow \infty$, and $\frac{k}{n} \rightarrow 0$, the error rate converges to B .

Exhaustive k-NN Algorithm:

- Compute squared distances from \tilde{q} to all n training points.
- Use a max-heap to maintain the k closest points.
- *Complexity:* Training: $\mathcal{O}(0)$, Query: $\mathcal{O}(nd + n \log k)$.

Voronoi Diagrams:

- The Voronoi cell for a point $w \in X$ is defined as:
 $\text{Vor}(w) = \{p \in \mathbb{R}^d : \|p - w\| \leq \|p - v\|, \forall v \in X\}$.
- Used for 1-NN to locate the nearest training point to a query point.
- *Complexity:* For low dimensions (2 - 5), efficient; for higher, consider k-d trees.

k-d Trees:

- Space-partitioning data structure to organize points in k-dimensional space.
- Splits along axis-aligned hyperplanes for efficient nearest neighbor search.
- *Building:* $\mathcal{O}(nd \log n)$, *Query:* sublinear, depends on dimensionality.

Key Advantages:

- Non-parametric and simple.
- Effective for small datasets and low dimensionality.

Challenges:

- High query time in large datasets.
- Curse of dimensionality: performance degrades as dimensionality increases.