

MODEL ERROR

Empirical Risk:

$$\hat{R}_{\mathcal{D}}(w) = \sum_{i=1}^n (w^\top x_i - y_i)^2$$

Generalisation Error (Pop. Risk):

$$L(f; \mathbb{P}_{X,Y}) = \mathbb{E}_{X,Y} \ell(f(X), Y)$$

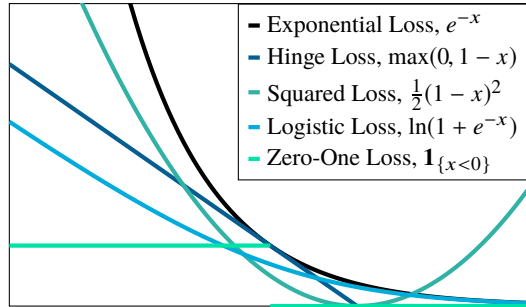
Bias-Variance Tradeoff:

$$\mathbb{E}_{\mathcal{D}}[L(\hat{f}; \cdot)] = \mathbb{E}_{X,Y}[(\hat{f}(X) - \mathbb{E}_{\mathcal{D}}[\hat{f}(X)])^2]$$

$$+ \mathbb{E}_X[(\mathbb{E}_{\mathcal{D}}[\hat{f}(X)] - f^*(X))^2] + \sigma$$

$$= \text{Var}_{\mathcal{D}}(\hat{f}) + \text{Bias}_{\mathcal{D}}^2(\hat{f}) + \text{Noise}$$

Least Squares: $X^\top X w = X^\top y$



REGULARIZATION

Lasso Regression (sparse):

$$\min_{w \in \mathbb{R}^d} (\|y - Xw\|_2^2 + \lambda \|w\|_1)$$

Ridge Regression (more precise):

$$\min_{w \in \mathbb{R}^d} (\|y - Xw\|_2^2 + \lambda \|w\|_2^2)$$

$$\nabla_w L(w) = 2X^\top(Xw - y) + 2\lambda w$$

Solution: $\hat{w} = (X^\top X + \lambda I)^{-1} X^\top y$

large $\lambda \Rightarrow$ larger bias, smaller variance

K-Fold Cross-Validation

Split Dataset into K sets (# methods), for each method, go through all sets and train it excluding that set and validating that set. Sum up all the validation errors of that method and choose smallest sum.

GRADIENT DESCENT

Converges only for convex case.

$$w^{t+1} = w^t - \eta_t \cdot \nabla \ell(w^t)$$

For linear regression:

$$\|w^t - w^*\|_2 \leq \|I - \eta X^\top X\|_{op}^t \|w^0 - w^*\|_2$$

$$\exists \eta \text{ with conv. speed } \rho = \|I - \eta X^\top X\|_{op}^t.$$

$$\eta_{\text{opt}} = \frac{2}{\lambda_{\min} + \lambda_{\max}} \text{ and max. } \eta \leq \frac{2}{\lambda_{\max}}.$$

Momentum:

$$w^{t+1} = w^t + \gamma \Delta w^{t-1} - \eta_t \nabla \ell(w^t)$$

MAXIMUM-MARGIN SOLUTION

If linearly separable, we can get:

$$w_{\text{MM}} := \arg \max_{\|w\|_2=1, w_0} \min_{1 \leq i \leq n} y_i (w^\top x_i + w_0)$$

Hard SVM

$$\hat{w} = \min_w \|w\|_2 \text{ s.t. } \forall_i, y_i w^\top x_i \geq 1$$

Soft SVM (allows "slack" in constraints)

$$\hat{w} = \min_{w, \xi} \frac{1}{2} \|w\|_2^2 + \lambda \sum_{i=1}^n \max(0, 1 - y_i w^\top x_i)$$

Metrics True Class γ err1/FPR: $\frac{\text{FP}}{\text{TN}+\text{FP}}$

$$\gamma = +1 \quad \gamma = -1 \quad \text{err2/FNR: } \frac{\text{FN}}{\text{TP}+\text{FN}}$$

$$f(x)=+1 \quad \text{TP} \quad \text{FP} \quad \text{Precision: } \frac{\text{TP}}{\text{TP}+\text{FP}}$$

$$f(x)=-1 \quad \text{FN} \quad \text{TN} \quad \text{TPR/Recall: } \frac{\text{TP}}{\text{TP}+\text{FN}}$$

ROC: Plot TPR=1-FNR vs. FPR and compare different ROC's with area under the curve.

F1-Score: $\frac{2\text{TP}}{2\text{TP}+\text{FP}+\text{FN}}$, Accuracy: $\frac{\text{TP}+\text{TN}}{\text{P}+\text{N}}$

Goal: large recall and small FPR.

KERNELS

$\exists \hat{\alpha}. \hat{w} = \Phi^\top \hat{\alpha}, K = \Phi \Phi^\top$, **Validity:**

1. K symmetric, $\forall x, z. k(x, z) = k(z, x)$

2. K positive semidef. (psd.), $\forall z. z^\top K z > 0$

lin.: $k(x, z) = x^\top z$, **poly.:** $k(x, z) = (x^\top z + 1)^m$

$$\exp\left(-\frac{\|x-z\|_2^2}{\tau}\right) = \begin{cases} \text{Laplacian Ker.} & p = 1 \\ \text{Gauss./RBF K.} & p = 2 \end{cases}$$

Composition Rules:

$$k = k_1 + k_2, \quad k = k_1 \cdot k_2, \quad c > 0 \Rightarrow k = c \cdot k_1$$

f convex; f polynomial or exp $\Rightarrow k = f(k_1)$

$$\forall f. k(x, y) = f(x)k_1(x, y)f(y)$$

Mercer's Theorem: Valid kernels can be decomposed into a lin. comb. of inner products.

Kern. Ridge Reg.: $\frac{1}{n} \|y - K\alpha\|_2^2 + \lambda \alpha^\top K \alpha$

KNN CLASSIFICATION

1. Pick k and distance metric d

2. For given x , find among $x_1, \dots, x_n \in D$ the k closest to $x \rightarrow x_{i_1}, \dots, x_{i_k}$

3. Output the majority vote of labels.

NEURAL NETWORKS

w are the weights and $\varphi: \mathbb{R} \mapsto \mathbb{R}$ is a *nonlinear activation function*: $\phi(x, w) = \varphi(w^\top x)$

ReLU: $\max(0, z)$, **Tanh:** $\frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$

Sigmoid: $\frac{1}{1 + \exp(-z)}$

Rand. init. weights by distr. assumption for φ .

ReLU: $2/n_{\text{in}}$; Tanh: $1/n_{\text{in}}$ or $1/(n_{\text{in}} + n_{\text{out}})$

Universal Approximation Theorem:

We can approximate any arbitrary smooth target function, with 1+ layer with sufficient width.

Forward Propagation

Input: $v^{(0)} = [x; 1]$ Output: $f = W^{(L)} v^{(L-1)}$

Hidden: $z^{(l)} = W^{(l)} v^{(l-1)}, v^{(l)} = [\varphi(z^{(l)}); 1]$

Backpropagation

Non-convex; **Reuse, Compute, Forward Pass**

$$(\nabla_{W^{(L)}} \ell)^\top = \frac{\partial \ell}{\partial W^{(L)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial W^{(L)}}$$

$$(\nabla_{W^{(L-1)}} \ell)^\top = \frac{\partial \ell}{\partial W^{(L-1)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial W^{(L-1)}}$$

$$(\nabla_{W^{(L-2)}} \ell)^\top = \frac{\partial \ell}{\partial W^{(L-2)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial z^{(L-2)}} \frac{\partial z^{(L-2)}}{\partial W^{(L-2)}}$$

Overfitting Prevention

Regularization: See lasso/ridge regression.

Early Stopping: Stops training upon converg.

Dropout: Deactiv. neurons rand. during train.

Batch Norm.: Norm. layer inputs $\mu = 0, \sigma = 1$.

CNNs $\varphi(W * v^{(l)})$

The output dimension when applying m different $f \times f$ filters to an $n \times n$ image with padding p and stride s is: $l = \frac{n+2p-f}{s} + 1$ For each channel there is a separate filter.

UNSUPERVISED LEARNING

k-Means Clustering

Optimization Goal (non-convex):

$$\hat{R}(\mu) = \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - \mu_j\|_2^2$$

Lloyd's heuristics: Init. cluster centers $\mu^{(0)}$

- Assign points to closest center
- Update μ_i as mean of assigned points

Converges in exp. time

Init. with k-Means++:

- Rand. data point $\mu_1 = x_i$
- Add μ_2, \dots, μ_k rand., with probability: Given $\mu_{1:j}$, pick $\mu_{j+1} = x_i$ where $p(i) = \frac{1}{z} \min_{l \in \{1, \dots, j\}} \|x_i - \mu_l\|_2^2$. E.g. further away from any centroid, higher chance.

Converges in expectation $O(\log k) \times \text{sol.}_{\text{opt}}$
Find k by negligible loss decrease or reg.

PRINCIPAL COMPONENT ANALYSIS

$$\arg \min_{W \in \mathbb{R}^{d \times k}: W^\top W = I_k} \sum_{i=1}^n \|x_i - W z_i\|_2^2,$$

where (v_i) are eigenvectors

$$W^* = (v_1 | \dots | v_k), \quad z_i^* = W^{*\top} x_i.$$

Principal eigenvector (v_1) shows into the direction of greatest variance in the data. The error is the deviation from it. Alternatively:

$$W^* = \arg \max_{W^\top W = I_k} \text{tr}(W^\top \Sigma W),$$

Where $\Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top$ is the empirical covariance. Closed form solution given by $w = v_1$ for $\lambda_1 \geq \dots \geq \lambda_d \geq 0: \Sigma = \sum_{i=1}^d \lambda_i v_i v_i^\top$. For $k > 1$, we only take the first k principal eigenvectors, such that $W^* = [v_1, \dots, v_k]$. In **SVD** the solution is given by the first k columns of V , with $X = U S V^\top$.

Kernel PCA

Ansatz: $w = \sum_{j=1}^n \alpha_j \phi(x_j) \Rightarrow \|w\|_2^2 = \alpha^\top K \alpha$

$$\alpha^* = \arg \max_{\alpha^\top K \alpha = 1} \alpha^\top K^\top K \alpha = \arg \max_{\alpha^\top \alpha = 1} \frac{\alpha^\top K^\top K \alpha}{\alpha^\top K \alpha}$$

Closed form, with $\lambda_1 \geq \dots \geq \lambda_n$:

$$\alpha^* = \frac{1}{\sqrt{\lambda_1}} v_1 \text{ for } K = \sum_{i=1}^n \lambda_i v_i v_i^\top$$

Autoencoders

Use a NN with smaller hidden layer than input size = output size to find a optimal subspace. $\min_x \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2$, lin. activ. f. \Rightarrow PCA

MLE & MAP

Maximum Likelihood Estimation (MLE)

(*) if discriminative, $p(a; b)$ since frequentist

$$\hat{\theta}_{\text{MLE}} := \arg \max_{\theta \in \Theta} p(\mathcal{D}; \theta)$$

$$\stackrel{\text{iid}}{=} \arg \max_{\theta \in \Theta} \prod_{i=1}^n p(x_i, y_i; \theta)$$

$$\stackrel{*}{=} \arg \min_{\theta \in \Theta} \sum_{i=1}^n -\log p(y_i | x_i; \theta),$$

Maximum A Posteriori Estimator (MAP)

(*) if discriminative, $p(a|b)$ since bayesian

$$\hat{\theta}_{\text{MAP}} := \arg \max_{\theta \in \Theta} p(\theta | \mathcal{D}) = \arg \max_{\theta \in \Theta} p(\mathcal{D} | \theta) p(\theta)$$

$$\stackrel{\text{iid}}{=} \arg \max_{\theta \in \Theta} (\prod_{i=1}^n p(x_i, y_i | \theta)) \cdot p(\theta)$$

$$\stackrel{*}{=} \arg \min_{\theta \in \Theta} \sum_{i=1}^n -\log p(y_i | x_i, \theta) - \log p(\theta)$$

Regression with MLE/MAP

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta \in \Theta} \sum_{i=1}^n (y_i - f(x_i; \theta))^2$$

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta \in \Theta} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 + \frac{\sigma_\epsilon^2}{\sigma_\theta^2} \|\theta\|_2^2$$

Regularization can be understood as MAP inference, with different priors (= regularizers) and likelihoods (= loss functions).

Statistical Models for Classification

f minimizing the population risk:

$$f^*(x) = \arg \max_{\hat{y}} p(\hat{y} | x)$$

This is called the Bayes' optimal predictor for the 0-1 loss. Assuming iid. Bernoulli noise, the conditional probability is:

$$p(y | \mathbf{x}, \mathbf{w}) \sim \text{Ber}(y; \sigma(\mathbf{w}^\top \mathbf{x}))$$

Where $\sigma(z) = \frac{1}{1+\exp(-z)}$ is the sigmoid function. Using MLE we get:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\text{argmin}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i))$$

Which is the logistic loss. Instead of MLE we can estimate MAP, e.g. with a Gaussian prior:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\text{argmin}} \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i})$$

BAYESIAN DECISION THEORY

Given $p(\mathbf{y} | \mathbf{x})$, a set of actions A and a cost $C : Y \times A \mapsto \mathbb{R}$, pick the action with the maximum expected utility.

$$a^* = \underset{a \in A}{\text{argmin}} \mathbb{E}_{\mathbf{y}}[C(\mathbf{y}, a) | \mathbf{x}]$$

Useful for asymmetric costs or abstention.

GENERATIVE MODELING (GM)

Aim to estimate $p(\mathbf{x}, \mathbf{y})$ for complex situations using Bayes' rule: $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x} | \mathbf{y}) \cdot p(\mathbf{y})$

Naive Bayes Model

GM for classification tasks. Assuming for a class label, each feature is independent. This helps estimating $p(\mathbf{x} | \mathbf{y}) = \prod_{i=1}^d p(x_i | y_i)$.

Gaussian Naive Bayes Classifier

Naive Bayes Model with Gaussians features. Estimate the parameters via MLE:

MLE for class prior: $p(\mathbf{y}) = \hat{p}_{\mathbf{y}} = \frac{\text{Count}(\mathbf{Y}=\mathbf{y})}{n}$

MLE for feature distribution:

Where: $p(x_i | \mathbf{y}) = \mathcal{N}(x_i; \hat{\mu}_{\mathbf{y},i}, \sigma_{\mathbf{y},i}^2)$

$$\mu_{\mathbf{y},i} = \frac{1}{\text{Count}(\mathbf{Y}=\mathbf{y})} \sum_{j|y_j=\mathbf{y}} x_{j,i}$$

$$\sigma_{\mathbf{y},i}^2 = \frac{1}{\text{Count}(\mathbf{Y}=\mathbf{y})} \sum_{j|y_j=\mathbf{y}} (x_{j,i} - \hat{\mu}_{\mathbf{y},i})^2$$

Predictions are made by:

$$\underset{\hat{\mathbf{y}}}{\text{argmax}} p(\hat{\mathbf{y}} | \mathbf{x}) = \underset{\hat{\mathbf{y}}}{\text{argmax}} p(\hat{\mathbf{y}}) \cdot \prod_{i=1}^d p(x_i | \hat{\mathbf{y}})$$

Equivalent to decision rule for bin. class.:

$$\mathbf{y} = \text{sgn} \left(\log \frac{p(\mathbf{Y} = +1 | \mathbf{x})}{p(\mathbf{Y} = -1 | \mathbf{x})} \right)$$

Where $f(\mathbf{x})$ is called the discriminant function. If the conditional independence assumption is violated, the classifier can be overconfident.

Gaussian Bayes Classifier

No independence assumption, model the features with a multivar. Gaussian $\mathcal{N}(\mathbf{x}; \mu_{\mathbf{y}}, \Sigma_{\mathbf{y}})$:

$$\mu_{\mathbf{y}} = \frac{1}{\text{Count}(\mathbf{Y}=\mathbf{y})} \sum_{j|y_j=\mathbf{y}} \mathbf{x}_j$$

$$\Sigma_{\mathbf{y}} = \frac{1}{\text{Count}(\mathbf{Y}=\mathbf{y})} \sum_{j|y_j=\mathbf{y}} (\mathbf{x}_j - \hat{\mu}_{\mathbf{y}})(\mathbf{x}_j - \hat{\mu}_{\mathbf{y}})^\top$$

This is also called the *quadratic discriminant analysis (QDA)*. LDA: $\Sigma_+ = \Sigma_-$, Fisher LDA:

$p(\mathbf{y}) = \frac{1}{2}$, Outlier detection: $p(\mathbf{x}) \leq \tau$.

Avoiding Overfitting

MLE is prone to overfitting. Avoid this by restricting model class (fewer parameters, e.g. GNB) or using priors (restrict param. values).

Discriminative models: $p(\mathbf{y} | \mathbf{x})$, fewer assumptions about data distribution

Generative models: $p(\mathbf{x}, \mathbf{y})$, can detect outliers, gen. missing data, less robust to outliers.

GAUSSIAN MIXTURE MODEL

Data is generated from a mixture of Gaussians:

$$p(\mathbf{x} | \theta) = \sum_{j=1}^k w_j \mathcal{N}(\mathbf{x}; \mu_j, \Sigma_j)$$

Estimate parameters by minimizing:

$$\underset{\theta}{\text{argmin}} - \sum_{i=1}^n \log \sum_{j=1}^k w_j \mathcal{N}(\mathbf{x}_i | \mu_j, \Sigma_j)$$

Non-convex objective, iteratively update parameters by predicting labels and imputing missing data.

Hard-EM Algorithm

E-Step: predict the most likely class for each data point:

$$z_i^{(t)} = \underset{z}{\text{argmax}} p(z | \mathbf{x}_i, \theta^{(t-1)})$$

$$= \underset{z}{\text{argmax}} p(z | \theta^{(t-1)}) \cdot p(\mathbf{x}_i | z, \theta^{(t-1)})$$

M-Step: compute MLE of $\theta^{(t)}$ as for GBC. *Problems:* Labels even if uncertain, tries to extract too much inf. Works poorly if clusters are overlapping. *Equivalent to k-Means with Lloyd's heuristics:* When having uniform weights and spherical covariances.

Soft-EM Algorithm

E-Step: Calculate the cluster membership weights for each point ($w_j = \pi_j = p(Z=j)$):

$$\gamma_j^{(t)}(\mathbf{x}_i) = p(Z=j | \mathbf{D}) = \frac{w_j \cdot p(\mathbf{x}_i; \theta_j^{(t-1)})}{\sum_k w_k \cdot p(\mathbf{x}_i; \theta_k^{(t-1)})}$$

M-Step: compute MLE with closed form:

$$w_j^{(t)} = \frac{1}{n} \sum_{i=1}^n \gamma_j^{(t)}(\mathbf{x}_i) \quad \mu_j^{(t)} = \frac{\sum_{i=1}^n \mathbf{x}_i \gamma_j^{(t)}(\mathbf{x}_i)}{\sum_{i=1}^n \gamma_j^{(t)}(\mathbf{x}_i)}$$

$$\Sigma_j^{(t)} = \frac{\sum_{i=1}^n \gamma_j^{(t)}(\mathbf{x}_i) (\mathbf{x}_i - \mu_j^{(t)}) (\mathbf{x}_i - \mu_j^{(t)})^\top}{\sum_{i=1}^n \gamma_j^{(t)}(\mathbf{x}_i)}$$

Init. the weights as uniformly distributed, rand. or with k-Means++ and for variances use spherical init. or empirical covariance of the data. Select k using cross-validation.

Special Cases of Gaussian Mixtures

• **Spherical:** Same variance in all directions.

$$\Sigma_k = \sigma_k^2 I; \text{ Parameters: } K$$

• **Diagonal:** Different variance for each dimension but no covariance.

$$\Sigma_k = \text{Diag}(\sigma_{k_1}^2, \sigma_{k_2}^2, \dots, \sigma_{k_d}^2); \text{ Par.: } K \cdot d$$

• **Tied:** Same cov. matrix for all components.

$$\Sigma_1 = \Sigma_2 = \dots = \Sigma_k; \text{ Parameters: } \frac{d(d+1)}{2}$$

• **Full:** Free covariance in all dimensions.

$$\omega, \mu, \Sigma; \text{ Parameters: } \frac{d(d+1)}{2} K$$

Degeneracy of GMMs

GMMs can overfit with limited data. To prevent this, add $v^2 I$ to the covariance matrices, preventing collapse (equivalent to using a Wishart prior). Choose v via cross-validation.

Gaussian-Mixture Bayes Classifiers

Assume that $p(x | y)$ for each class can be modelled by a GMM.

$$p(\mathbf{x} | y) = \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(\mathbf{x}; \mu_j^{(y)}, \Sigma_j^{(y)})$$

Giving highly complex decision boundaries:

$$p(y | \mathbf{x}) = \frac{1}{z} p(y) \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(\mathbf{x}; \mu_j^{(y)}, \Sigma_j^{(y)}).$$

GMMs for Density Estimation

Can be used for anomaly detection or data imputation. Detect outliers, by comparing the estimated density against τ . Allows to control the FP rate. Use ROC curve as evaluation criterion and optimize using CV to find τ .

General EM Algorithm

E-Step: Take the expected value over latent variables z to generate likelihood function Q : $Q(\theta; \theta^{(t-1)}) = \mathbb{E}_Z[\log p(\mathbf{X}, \mathbf{Z} | \theta) | \mathbf{X}, \theta^{(t-1)}]$

$$= \sum_{i=1}^n \sum_{z_i=1}^k \gamma_{z_i}(\mathbf{x}_i) \log p(\mathbf{x}_i, z_i | \theta)$$

with $\gamma_z(\mathbf{x}) = p(z | \mathbf{x}, \theta^{(t-1)})$.

M-Step: Compute MLE / Maximize:

$$\theta^{(t)} = \underset{\theta}{\text{argmax}} Q(\theta; \theta^{(t-1)})$$

We have monotonic convergence, each EM-iteration increases the data likelihood.

GAN

New loss: $\min_{w_G} \max_{w_D} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D(\mathbf{x}, w_D)] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D(G(\mathbf{z}, w_G), w_D))]$

- **Saddle Point:** Training seeks a saddle point.
- **Capacity:** Conv. if \mathbf{G} and \mathbf{D} have enough capacity.
- **Optimal \mathbf{D} for fixed \mathbf{G} :**

$$D_G(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})}$$

- **Fake Probability:** $1 - D_G$.
- **Issues:** Oscill., divergence, mode collapse.

Performance Metric:

$$DG = \max_{w_D} M(w_G, w_D') - \min_{w_G'} M(w_G', w_D)$$

where $M(w_G, w_D)$ is the training objective. ““

VARIOUS

Derivatives:

$$\nabla_{\mathbf{x}} \mathbf{x}^\top \mathbf{A} = \mathbf{A} \quad \nabla_{\mathbf{x}} \mathbf{a}^\top \mathbf{x} = \nabla_{\mathbf{x}} \mathbf{x}^\top \mathbf{a} = \mathbf{a}$$

$$\nabla_{\mathbf{x}} \mathbf{b}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}, \quad \nabla_{\mathbf{x}} \mathbf{x}^\top \mathbf{x} = 2\mathbf{x}, \quad \nabla_{\mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$$

$$\text{Square Loss: } \nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X} \mathbf{w}\|_2^2 = 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y}$$

Bayes Theorem:

$$p(y | x) = \frac{1}{p(x)} p(y) \cdot p(x | y)$$

Normal Distribution:

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{(\mathbf{x}-\mu)^\top \Sigma^{-1}(\mathbf{x}-\mu)}{2}\right)$$

Exponential Distribution: $\text{Exp}(\lambda) = \lambda e^{-\lambda x}$

Other Facts

Memoryless: $p(X > a + b | X \geq a) = p(X > b)$

Tower Property: $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X | Y]]$

$$\Rightarrow \mathbb{E}[X] = \sum_{y \in \mathcal{Y}} \mathbb{E}[X | Y = y] p_Y(y)$$

$$\text{Tr}(AB) = \text{Tr}(BA), \quad \text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2,$$

$$\mathbf{X} \in \mathbb{R}^{n \times d} : \mathbf{X}^{-1} \rightarrow O(d^3); \mathbf{X}^\top \mathbf{X} \rightarrow O(nd^2),$$

$$\binom{n}{k} = \frac{n!}{(n-k)!k!}, \quad \|\mathbf{w}^\top \mathbf{w}\|_2 = \sqrt{\mathbf{w}^\top \mathbf{w}}$$

$$\text{Cov}[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top]$$

$$p(\mathbf{z} | \mathbf{x}, \theta) = \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{p(\mathbf{x} | \theta)}$$

Convexity

- $\alpha f + \beta g, \alpha, \beta \geq 0$, convex if f, g convex
- $f \circ g$, convex if $[f$ convex and g affine (e.g. $ax + b$)] or $[f$ non-decreasing and g convex]
- $\max(f, g)$, convex if f, g convex
- $L(\lambda \mathbf{w} + (1 - \lambda) \mathbf{v}) \leq \lambda L(\mathbf{w}) + (1 - \lambda) L(\mathbf{v})$

1. Order: $L(\mathbf{w}) + \nabla L(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) \leq L(\mathbf{v})$

2. Order: Hessian $\nabla^2 L(\mathbf{w}) \geq 0$ (psd)