

Math

Definition (Positive Semidefinite, PSD):

M is PSD if $\langle v, Mv \rangle = v^\top M v \geq 0, \forall v \in \mathbb{R}^d, M \in \mathbb{R}^{d \times d}$

Proposition (Convex Functions): A twice differentiable function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if the Hessian $\nabla^2 f(x)$ is positive semi-definite for all $x \in \mathbb{R}^n$

Proposition (Fermat's necessary condition): If w is a local extreme of f , then $\nabla f(w) = 0$. If f is convex the condition is \Leftrightarrow .

Halfspaces & The Perceptron Algorithm

Definition (Binary classification): Given n known pairs $\{(\mathbf{x}_i, y_i)\}$ where $x_i \in \mathbb{R}^d, y_i \in \{\pm 1\}$ we want to learn a classification rule

$h(\mathbf{x}) = y$ Such that,
 $\mathbb{P}_{(\mathbf{x}, y) \sim \rho}[h(\mathbf{x}) = y]$ is Large

Definition (Affine Set/Hyper Planes): The signed distance of a vector \mathbf{x} is given by

$h(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$
 \mathbf{w} : Feature Weights, b : Bias

Algorithm (Perceptron):

```
def perceptron(D, d, delta, k):
    w = np.zeros(d) # Initialize w
    b = 0 # Initialize b
    for t in range(k):
        # Receive training example (xt, yt)
        xt, yt = D[t]
        # Compute prediction ŷ
        y_hat = np.dot(w, xt) + b
        # Check the update condition
        if yt * y_hat <= delta:
            # Update w and b
            w = w + yt * xt
            b = b + yt
    return w, b
```

Note (Padding): $\vec{x}_i' = (x_1, x_2, \dots, x_d, 1)^\top$,
 $\mathbf{w} = (w_1, w_2, \dots, w_d, b)$. We can write $\langle \mathbf{w}, \vec{x}' \rangle$

Definition (Linear separability): D is linearly separable if $\exists \vec{u}, \|\vec{u}\| = 1$, and a margin $\gamma > 0$ such that

$$\forall (\mathbf{x}, y) \in D, y(\vec{u}^\top \mathbf{x}) > \gamma$$

Theorem (Perceptron Convergence theorem): If D is linearly separable with margin $\gamma > 0$ and $\forall (\vec{x}, y) \in D, \|\vec{x}\| \leq R$ the Perceptron algorithm will converge by update

$$k^2 \gamma^2 \leq \|\mathbf{w}_k\|^2 \leq k R^2,$$

$$k \leq \frac{R^2}{\gamma^2}$$

Note (Multiclass):

- one-vs-all: Train one for each class output $\arg \max_{k=1 \dots c} (\mathbf{w}_k^\top \mathbf{x} + b_k)$
- one-vs-one: For each pair (k, k') train $\mathbf{w}_{k,k'}$ total of $\binom{c}{2}$ perceptrons, output majority vote.

Example (Loss Functions):

$\max(0, 1 - yf(\mathbf{x})) = (1 - yf(\mathbf{x}))^+$ (Hinge Loss)
 $\log(1 + e^{-yf(\mathbf{x})})$ (Cross Entropy Loss)
 $(y - f(\mathbf{x}))^2$ (Squared Error)

Support Vector Machines

Definition (Maximum Margin Classifier): For dataset $\mathcal{D} = \{(\vec{x}_i, y_i)\}$, $y_i \in \{\pm 1\}$ linearly separated by a hyperplane $w^\top x = -b$ with margin γ

$$\gamma = \underbrace{\min_i \frac{1}{\|\mathbf{w}\|} y_i (w^\top x_i + b)}_{\text{Min Distance from the hyperplane}} = \frac{1}{\|\mathbf{w}\|} \min_i y_i (w^\top x_i + b)$$

The maximum margin classifier problem is defined as

$$\hat{w}, \hat{b} = \arg \max_{w, b} \frac{1}{\|\mathbf{w}\|} \min_i y_i (w^\top x_i + b)$$

Since scaling w, b does not affect the distance we have:

$$\hat{w}, \hat{b} = \arg \max_{w, b} \frac{1}{\|\mathbf{w}\|} \text{ subject to: } y_i (w^\top x_i + b) \geq 1 \forall i$$

$$\hat{w}, \hat{b} = \arg \min_{w, b} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to: } y_i (w^\top x_i + b) \geq 1 \forall i$$

Note (Lagrangian dual of maximum margin):

$$\min_{w, b} \max_{\lambda \geq 0} \mathcal{L}(w, b, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \lambda_i (y_i \hat{y}_i - 1)$$

Let $\hat{y}_i = (w^\top x_i + b)$

KKT Conditions:

- Primal feasibility: $y_i \hat{y}_i - 1 \geq 0$ (Each datapoint is correctly classified)
- Dual feasibility: $\lambda_i \geq 0$
- Complementary slackness: $\lambda_i (y_i \hat{y}_i - 1) \geq 0$, either $\lambda_i = 0$ (inactive), when $\lambda_i > 0$ then its directly influencing the decision boundary.

Note (Hard vs Soft Margin SVM):

Hard	Soft
$\arg \min_{w, b} \frac{1}{2} \ \mathbf{w}\ ^2$	$\arg \min_{w, b} \frac{1}{2} \ \mathbf{w}\ ^2 + C \sum_i (y_i \hat{y}_i - 1)^+$
Subject to $y_i \hat{y}_i \geq 1$	Hinge Loss

$$\min_{\lambda_i \geq 0} - \sum_i \lambda_i + \min_{0 \leq \lambda \leq C} - \sum_i \lambda_i +$$

$$\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \quad \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle$$

$$\mathbf{w}^* = \sum_{i=1}^n \lambda_i y_i x_i, \quad \lambda_i \geq 0 \quad \mathbf{w}^* = \sum_{i=1}^n \lambda_i y_i x_i, \quad 0 \leq \lambda_i \leq C$$

$$\mathbf{b}^* = \frac{1}{N_{SV}} \sum_{i \in SV} (y_i - \mathbf{w}^\top x) \quad \mathbf{b}^* = \frac{1}{N_{SV}} \sum_{i \in SV} (y_i - \mathbf{w}^\top x)$$

High Variance $y_i \hat{y}_i < 0$: Incorrectly classified,
 $0 \leq y_i \hat{y}_i < 1$: Weakly correct,
 $y_i \hat{y}_i > 1$: Strong correct,

Kernel Methods

Theorem (Mercer kernels):

- Any symmetric function $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a kernel if and only if there exists some $\phi: \mathcal{X} \mapsto \mathcal{H}$ such that $k(x, x') = \langle \phi(x), \phi(x') \rangle$
- A function $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a Mercer kernel \Leftrightarrow for any $n \in \mathbb{N}, x_i \in \mathcal{X}$ the kernel matrix \mathbf{K} with $K_{ij} = k(x_i, x_j)$ is Symmetric and PSD

Example (Kernels):

$\exp(-\gamma \|x - x'\|_2^2)$ (Gaussian)
 $\exp(-\gamma \|x - x'\|)$ (Laplace)

Note (Kernel Trick SVM):

$$\mathbf{b}^* = \frac{1}{N_{SV}} \sum_{x_i \in SV} \left(y_i - \sum_{x_j \in SV} \lambda_j y_j K_{ij} \right), \quad K_{ij} = k(x_i, x_j)$$

$$\hat{y} = \sum_{x_i \in SV} \lambda_i y_i k(x, x_i) + \mathbf{b}^*$$

Linear Regression

Definition (Regression): Given n known pairs $\{(\mathbf{x}_i, y_i)\}$ where $x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$ and $(\vec{x}_i, y_i) \stackrel{\text{iid}}{\sim} P$, we want to learn $h: \mathbb{R}^d \rightarrow \mathbb{R}$

$\mathbb{E}_{(\vec{x}, y) \sim P}[l_w(\vec{x}, y)]$ is small

$$\Rightarrow \arg \min_w \frac{1}{n} \sum l_w(\vec{x}_i, y)$$

Where $\ell_w(\cdot, \cdot)$ is the loss function

Note (Calculus): Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla f(\vec{x}): \mathbb{R}^d = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right) \text{ (Gradient)}$$

$$\nabla^2 f(\vec{x}): \mathbb{R}^{d \times d} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_d^2} \end{pmatrix}$$

Note (Equivalent notation of loss): Let $A \in \mathbb{R}^{n \times (d+1)}$ be a matrix of padded feature vectors. Let $z \in \mathbb{R}^{n \times 1}$ be the vector of outputs in training set. Then we can write total loss as

$$L = \sum_{i=1}^n (\langle \mathbf{w}, \vec{x}' \rangle - y_i)^2$$

$$L = \|\mathbf{A} \mathbf{w} - \mathbf{z}\|_2^2$$

Note (least-squares loss):

$$L = \|\mathbf{X} \mathbf{w} - \vec{y}\|_2^2$$

$$\nabla_w L = 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \vec{y}$$

$$\nabla_w^2 L = 2\mathbf{X}^\top \mathbf{X} \Rightarrow \nabla^2 \mathbf{X}^\top \mathbf{X} v = 2\|\mathbf{X} v\|_2^2 \geq 0$$

least-squares loss convex. Solving for $\nabla_w L = 0$

$$2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \vec{y} = \vec{0} \rightarrow \hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{y}$$

Note (Regularization):

L_1 Lasso	L_2 Ridge/Tikhonov
$\min_{\mathbf{w}} \ \mathbf{A} \mathbf{w} - \vec{y}\ _2^2 + \lambda \ \mathbf{w}\ _1$	$\min_{\mathbf{w}} \ \mathbf{A} \mathbf{w} - \vec{y}\ _2^2 + \lambda \ \mathbf{w}\ _2^2$
<ul style="list-style-type: none"> Penalizes non-zero weights. Many w_i will be 0. Weights \mathbf{w} will be sparse. 	<ul style="list-style-type: none"> Penalizes large weights, stabilizes the weights.
no closed-form	$\hat{\mathbf{w}}(\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}) = \mathbf{X}^\top \mathbf{Y}$

Definition (Entropy):

$$h(X) = \mathbb{E}(-\log p(X))$$

$$= \sum_{i=1}^n \Pr[x_i] \log \Pr[x_i] = \int_{\mathcal{X}} p(x) \log p(x) dx$$

Note (Kullback-Leibler divergence):

$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}$$

$$D_{KL}(P \parallel Q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

measures dissimilarity between a reference and model distribution (relative entropy)

Logistic regression

Note (Sigmoid Function):

$$h(x) = \frac{1}{1 + \exp(-\langle \mathbf{w}, \mathbf{x} \rangle)} = \text{sigmoid}(\langle \mathbf{w}, \mathbf{x} \rangle)$$

Note (Logistic Regression): log-likelihood function for logistic regression is:

$$\log \mathcal{L}(w \mid \mathbf{x}, \mathbf{y}) = \sum_{i=1}^n -\log(1 + e^{-\tilde{y}_i \langle \mathbf{w}, \mathbf{x}_i \rangle})$$

$$\tilde{y}_i = \begin{cases} -1 & \text{if } y_i = 0 \\ +1 & \text{if } y_i = 1 \end{cases}$$

Cross-entropy loss $l_w(\mathbf{x}, \tilde{y}) = \log(1 + e^{-\tilde{y} \langle \mathbf{w}, \mathbf{x} \rangle})$

Definition (Gradient descent): In Gradient descent the parameter vector w is updated in the direction opposite to ∇f with step size η (learning rate)

$$w_t = w_{t-1} - \eta \nabla f(w_{t-1})$$

For a loss function this is computed over a batch of n

$$\nabla l_w(\mathbf{x}, y) = \frac{1}{n} \sum_{i=1}^n \nabla l_{w_{t-1}}(\mathbf{x}_i, y_i)$$

In Stochastic gradient descent the gradient is estimated using a random sample (minibatch) of size m

$$w_t = w_{t-1} - \underbrace{(\nabla^2 l_{w_{t-1}})^{-1}}_{\text{Inverse Hessian}} \nabla l_{w_{t-1}} \quad \text{(Newton's method)}$$

Note (Linear vs Logistic):

	Linear	Logistic
• Loss	$\sum_{i=1}^n \underbrace{(y_i - \hat{y}_i)^2}_{\text{Square Error}}$	$\sum_{i=1}^n \underbrace{\log(1 + e^{-\tilde{y}_i \langle \mathbf{w}, \mathbf{x}_i \rangle})}_{\text{Cross Entropy}}$
• Predict	$\hat{y}_i = \langle \mathbf{w}, \mathbf{x}_i \rangle$	$\hat{p}_i = \text{sigmoid}(\langle \mathbf{w}, \mathbf{x}_i \rangle)$
• Objective	$\ \mathbf{y} - \hat{\mathbf{y}}\ _2^2$	$\text{KL}(\frac{1+\mathbf{y}}{2} \parallel \hat{\mathbf{p}})$
• Gradient	$\hat{\mathbf{w}} - \eta \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})$	$\hat{\mathbf{w}} - \eta \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})$

Note (Multiclass Logistic Regression):

$$\Pr[y = k \mid \mathbf{x}, \mathbf{w}] = \frac{e^{\langle \mathbf{w}_k, \mathbf{x} \rangle}}{\sum_{i=1}^c e^{\langle \mathbf{w}_i, \mathbf{x} \rangle}} = \text{softmax}$$

Non-Parametric Methods

Definition (Kerne Density Estimation): Given $\{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}^d$ the Parzen estimate for probability at point x_o is

$$\hat{p}(\mathbf{x}_o) = \frac{1}{n\lambda} \sum_{i=1}^n K_\lambda(\mathbf{x}_o, x_i)$$

Where $K_\lambda: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^+$ is the kernel function. λ is a hyperparameter called length scale (or bandwidth). K_λ should be symmetric and shift invariant.

$$K_{\lambda(x_o, x_i)} = K\left(\frac{x_o}{\lambda} - \frac{x_i}{\lambda}\right)$$

k-Nearest Neighbors

Algorithm (KNN):

```
import numpy as np
def k_nearest_neighbors(data, k, x):
    distances = []
    for xi, yi in data:
        # Euclidean norm
        distance = np.linalg.norm(x - xi)
        distances += [(distance, yi)]
        distances.sort(key=lambda x: x[0])
    k_nearest_labels = distances[:k]
    majority_label = count(k_nearest_labels)
    return majority_label
```

Training: time: $O(1)$, space: $O(nd)$

Testing: time $O(nd + n \log k)$, space: $O(nd)$

Note (An upper bound on KNN error rate): Cover-Hart Theorem, for classifier with $k = 1$

$$\lim_{n \rightarrow \infty} \epsilon_{1\text{-NN}} < 2\epsilon_{\text{Bayes}}(1 - \epsilon_{\text{Bayes}})$$

- Number of samples needed to train KNN to reach lower bound on error grows exponentially. Euclidean distance is unhelpful in high dimensions all vectors are almost equidistant

K-Means Clustering

Algorithm (K-Means):

```
def kmeans(X, k, max_iter=300):
    N = len(X)
    C0 = randint(0, high=k, size=N)
    C1 = zeros(N)
    iter = 0
    while iter < max_iter:
        # Centroids of each cluster
        mu = [mean(X[CO==i]) for i in 0..k]
        for i in range(N):
            # Best centroid for x_i
            best = argmin([norm(X[i] - mu[j])**2 for j in 0..k])
            C1[i] = best
        if array_equal(C0, C1):
            # No updates detected
            break
        C0 = C1
        iter += 1
    return C0
```

Bayesian Learning

Definition (Prior and Posterior): In a Bayesian approach assume the parameter θ follows a prior pdf $\Pr(\theta)$.

Given a prior pdf $\Pr(\theta)$, after observing some data \mathcal{D} belief on the probable values of θ will have changed. We get the posterior:

$$\Pr[\theta|\mathcal{D}] = \frac{\Pr[\mathcal{D}|\theta] \Pr[\theta]}{\Pr[\mathcal{D}]} = \frac{\Pr[\mathcal{D}|\theta] \Pr[\theta]}{\int \Pr[\mathcal{D}|\theta] \Pr[\theta] d\theta}$$

Note: computing the integral may not be tractable.

Theorem (Bayes classifier): A classification problem with $\vec{x} \in \mathbb{R}^d$ and $Y \in [c] := \{1 \dots c\}$ The optimal (Bayes) classifier is

$$\begin{aligned} \underset{h: \mathbb{R}^d \rightarrow [c]}{\operatorname{argmin}} \Pr(Y \neq h(\vec{x})) \\ h^*(\vec{x}) = \underset{k \in [c]}{\operatorname{argmax}} \Pr(Y = k \mid \vec{x}) \\ \underset{k \in [c]}{\operatorname{argmax}} \underbrace{\Pr(X = \vec{x} \mid Y = k)}_{\text{likelihood}} \underbrace{\Pr(Y = k)}_{\text{prior}} \end{aligned}$$

Definition (Maximum likelihood estimate (MLE)): MLE estimate maximizes the likelihood of observing the data \mathcal{D} . This ignores any prior on θ .

$$\begin{aligned} L(\theta; \mathcal{D}) &:= \Pr(\mathcal{D}|\theta) = \prod_{i=1}^n \Pr(\mathbf{x}_i|\theta) \\ \theta_{\text{MLE}} &= \underset{\theta \in \Theta}{\operatorname{argmax}} L(\theta; \mathcal{D}) \\ \theta_{\text{MLE}} &= \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^n -\log \Pr(\mathbf{x}_i|\theta) \end{aligned}$$

Definition (Maximum a posteriori (MAP)):

$$\begin{aligned} \theta_{\text{MAP}} &= \underset{\theta \in \Theta}{\operatorname{argmax}} \Pr(\theta \mid \mathcal{D}) \\ \rightarrow \underset{\theta \in \Theta}{\operatorname{argmax}} \underbrace{-\log \Pr(\mathcal{D}|\theta)}_{\text{negative log-likelihood}} - \underbrace{\log \Pr(\theta)}_{\text{prior as regularization}} \end{aligned}$$

Example (Ridge Regression and MAP): Linear regression with λ is equivalent to MAP estimate of \mathbf{w} using a prior on $\mathbf{w} \sim \mathcal{N}_d(\mathbf{0}, \Sigma)$ where $\Sigma^{-1} = \lambda I$

Example (Bayesian linear regression):

$$\begin{aligned} \Pr(\mathbf{w}) &= \mathcal{N}(\vec{0}, \Sigma) \quad (\text{Gaussian prior on the weights}) \\ \Pr(\epsilon) &= \mathcal{N}(0, \sigma^2) \\ \Pr(\mathbf{w}|\mathcal{D}) &\propto \Pr(\mathbf{w}) \Pr(\mathcal{D}|\mathbf{w}) \propto \mathcal{N}(\bar{\mathbf{w}}, A^{-1}) \quad (\text{Posterior}) \\ A &= \frac{1}{\sigma^2} X^\top X + \Sigma^{-1}, \quad \bar{\mathbf{w}} = A^{-1} \frac{1}{\sigma^2} X^\top \mathbf{y} \\ \Pr(y^*|\vec{x}^*, \mathcal{D}) &= \int_{\mathbf{w}} \Pr(y^*|\vec{x}^*, \mathbf{w}) \Pr(\mathbf{w}|\mathcal{D}) \quad (\text{Prediction}) \\ \Pr(y^*|\vec{x}^*, \mathcal{D}) &= \mathcal{N}(\langle \vec{x}^*, \bar{\mathbf{w}} \rangle, \sigma^2 + \langle \vec{x}^*, A^{-1} \vec{x}^* \rangle) \end{aligned}$$

Gaussian Processes

Definition (Gaussian Processes): Function space view $\Pr(f|\mathcal{D})$ posterior over the function f (Infinite dimensional Gaussian Distribution).

$$f(\mathbf{x}) \sim \text{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \forall \mathbf{x}', \mathbf{x}$$

Gaussian Processes are specified completely by μ, Σ

$$\begin{aligned} m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})], \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \end{aligned}$$

k is the kernel covariance function.

Example (Gaussian Processes): Let $\mathbf{y}(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$, consider a prior on $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$

$$\begin{aligned} \mathbf{y} &= \Phi \mathbf{w} \\ \mathbb{E}[\mathbf{y}] &= \Phi E[\mathbf{w}] = \mathbf{0} \\ \text{Cov}[\mathbf{y}] &= \mathbb{E}[\mathbf{y}^\top \mathbf{y}] = \frac{1}{\alpha} \Phi \Phi^\top = \mathbf{K} \end{aligned}$$

Where \mathbf{K} is the Gram matrix with

$$\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\alpha} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

Note (Gaussian processes for regression):

$$\underset{\text{True Value}}{t_n} = \underset{\mathbf{y}(\vec{x}_n)}{y_n} + \underset{\text{Noise}}{\epsilon_n}$$

Assuming the noise process is Gaussian

$$\Pr[\mathbf{t}|\mathbf{y}] = \mathcal{N}(\mathbf{t}|\mathbf{y}, \beta^{-1} \mathbf{I}_N)$$

By definition of the gaussian process marginal pdf of \mathbf{y} is defined by a Gram matrix \mathbf{K}

$$\begin{aligned} \Pr[\mathbf{y}] &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \alpha^{-1} \mathbf{K}) \\ \Pr[\mathbf{t}] &= \int \Pr(\mathbf{t}|\mathbf{y}) \Pr[\mathbf{y}] d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \\ \mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) &= \alpha^{-1} k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1} \times \mathbb{1}[n = m] \end{aligned}$$

For prediction we augment the covariance matrix

$$\begin{aligned} \Pr[\mathbf{t}_{N+1}] &= \mathcal{N}(\mathbf{t}_{N+1} \mid \mathbf{0}, \mathbf{C}_{N+1}) \\ \mathbf{C}_{N+1} &= \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^\top & c \end{pmatrix}, \quad \mathbf{k} \in \mathbb{R}^N, \quad \mathbf{k}_i = \alpha^{-1} k(\mathbf{x}_{N+1}, \mathbf{x}_i) \\ c &= \alpha^{-1} k(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) + \beta^{-1} \end{aligned}$$

Then the predictive distribution is

$$\begin{aligned} \Pr[t_{N+1}|\mathbf{t}] &= \mathcal{N}(\mu(\mathbf{x}_{N+1}), \sigma^2(\mathbf{x}_{N+1})) \\ \mu(\mathbf{x}_{N+1}) &= \mathbf{k}^\top \underbrace{\mathbf{C}_N^{-1} \mathbf{t}}_{\vec{a}}, \quad \sigma^2(\mathbf{x}_{N+1}) = c - \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{k} \\ \mu(\mathbf{x}_{N+1}) &= \langle \mathbf{k}, \vec{a} \rangle = \sum_{i=1}^N k(\mathbf{x}_i, \mathbf{x}_{N+1}) \vec{a}_i \end{aligned}$$

Note (Gaussian Process Time Complexity):

Type	Fixed Basis Functions m Features	Gaussian Process N Data Points
Train	$O(M^3)$	$O(N^3)$
Test	$O(M^2)$	$O(N^2)$

Decision Trees

Definition (Decision Trees Growing): Growing a tree means we choose a dimension or feature j^* and a threshold t^* to split that minimized the loss ℓ .

$$\begin{aligned} L_{j,t} &= \{(\mathbf{x}_i, y_i) \in S : \mathbf{x}_{ij} \leq t\} \\ R_{j,t} &= \{(\mathbf{x}_i, y_i) \in S : \mathbf{x}_{ij} > t\}, \\ (j^*, t^*) &= \underset{j=1..d}{\operatorname{argmin}} \min_{t \in T_j} |L_{j,t}| \ell(L_{j,t}) + |R_{j,t}| \ell(R_{j,t}) \end{aligned}$$

Note (Decision Trees Classification Cost): For a region \mathcal{D} we define

$$\begin{aligned} \hat{p}_k &= \frac{1}{|\mathcal{D}|} \sum_{y_i \in \mathcal{D}} \mathbb{1}(y_i \in k) \quad (\text{Empirical fraction of class } k) \\ \hat{y}_i &= \underset{k=1..c}{\operatorname{argmax}} \hat{p}_k \quad (\text{Majority class in } \mathcal{D}) \end{aligned}$$

Loss(\mathcal{D}) = $1 - \max_k \hat{p}_k = 1 - \hat{p}_y$ (Misclassification error)

$$\text{Loss}(\mathcal{D}) = \sum_{k=1}^c \hat{p}_k (1 - \hat{p}_k) \quad (\text{Gini index})$$

$$\text{Loss}(\mathcal{D}) = - \sum_{k=1}^c \hat{p}_k \log(\hat{p}_k) \quad (\text{Entropy})$$

Note (Minimal Cost-Complexity Pruning): Grow a tree fully first then regularize using a hyperparameter α .

$$R_\alpha(T) = R(T) + \alpha |\tilde{T}|$$

Where $|\tilde{T}|$ is the number of terminal nodes in T and $R(T)$ is the total loss for the leaves. Penalizes deep trees/more leaves.