

Advanced Machine Learning

Cristian Perez Jensen

February 7, 2025

Note that these are not the official lecture notes of the course, but only notes written by a student of the course. As such, there might be mistakes. The source code can be found at github.com/cristianpjensen/eth-cs-notes. If you find a mistake, please create an issue or open a pull request.

Contents

1	Paradigms of data science	1
2	Anomaly detection	2
2.1	Principal component analysis	2
2.2	Gaussian mixture model	3
2.3	Summary	5
3	Density estimation	6
4	Regression	9
4.1	Linear regression	9
4.2	Regularization	11
4.3	Polynomial regression	11
5	Causality	13
5.1	Counterfactual invariance	13
6	Gaussian processes	15
6.1	Kernels	16
6.2	Prediction	17
7	Uncertainty quantification	18
7.1	Statistical model validation	18
7.2	Uncertainty in linear models	20
7.3	Statistical testing	20
7.4	Bayesian neural networks	21
7.5	Information-based transductive learning	23
8	Convex optimization	27
8.1	Support vector machine	29
8.2	Soft-margin SVM	32
8.3	Kernelization	32
8.4	Multi-class SVM	33
8.5	Structural SVM	33
9	Ensembles	35
9.1	Bagging	36
9.2	Random forests	36
9.3	AdaBoost	36
10	Stable diffusion	39
10.1	Diffusion models	39
10.2	U-net	40
10.3	Latent diffusion models	40
10.4	Text embeddings	41
10.5	Cross-attention	41
11	Non-parametric Bayesian methods	42
11.1	Dirichlet process	42
11.2	Chinese restaurant process	43

11.3	<i>Dirichlet process mixture model</i>	44
12	<i>Statistical learning theory</i>	46
12.1	<i>PAC learning</i>	46
12.2	<i>Axis-aligned rectangles</i>	47
12.3	<i>Stochastic setting</i>	48
12.4	<i>Vapnik-Chervonenkis</i>	48

List of symbols

\doteq	Equality by definition
\approx	Approximate equality
\propto	Proportional to (here, generally used to remove any terms that do not matter for an argmin or argmax operator)
\mathbb{N}	Set of natural numbers
\mathbb{R}	Set of real numbers
$i : j$	Set of natural numbers between i and j
$f : A \rightarrow B$	Function f that maps elements of set A to elements of set B
$\mathbb{1}\{\text{predicate}\}$	Indicator function (1 if predicate is true, otherwise 0)
$\mathbf{v} \in \mathbb{R}^n$	n -dimensional vector
$\mathbf{M} \in \mathbb{R}^{m \times n}$	$m \times n$ matrix
\mathbf{M}^\top	Transpose of matrix \mathbf{M}
\mathbf{M}^{-1}	Inverse of matrix \mathbf{M}
$\det(\mathbf{M})$	Determinant of \mathbf{M}
$\frac{d}{dx}f(x)$	Ordinary derivative of $f(x)$ w.r.t. x at point $x \in \mathbb{R}$
$\frac{\partial}{\partial x}f(\mathbf{x})$	Partial derivative of $f(\mathbf{x})$ w.r.t. x at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}f(\mathbf{x}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$
$\nabla_{\mathbf{x}}^2f(\mathbf{x}) \in \mathbb{R}^{n \times n}$	Hessian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point $\mathbf{x} \in \mathbb{R}^n$

1 Paradigms of data science

Let $\{x_1, \dots, x_n\}$ be *i.i.d.* samples, generated by an unknown distribution P . Assume that this distribution is in a distribution family,

$$\mathcal{H} = \{p(\cdot \mid \theta) \mid \theta \in \Theta\}.$$

The goal is to learn the parameters θ that fit the data $\{x_1, \dots, x_n\}$ best.

Frequentism. In frequentism, the MLE (*Maximum Likelihood Estimator*) parameters maximize the following,

$$\begin{aligned} \theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\} \mid \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log p(x_i \mid \theta). \end{aligned}$$

Bayesianism. Bayesianism assumes that there is a prior over distributions. The MAP (*Maximum A Posteriori*) parameters maximize the following,

$$\begin{aligned} \theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\theta \mid \mathbf{X}) \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\} \mid \theta) \cdot p(\theta) && \text{Bayes' rule.} \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\theta) + \sum_{i=1}^n \log p(x_i \mid \theta). \end{aligned}$$

In practice, the prior acts as a regularization term, such that the posterior does not move too far from the prior.

Statistical learning. Now, assume that we have labeled samples $(x, y) \in \mathcal{X} \times \mathcal{Y}$, where y is the target variable. Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. For a predictor function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we define its risk as the expected loss,

$$\mathcal{R}(f) \doteq \mathbb{E}_{X,Y}[\ell(y, f(x))].$$

In statistical learning, we want to find a function that minimizes the risk. However, since the distribution over X, Y is unknown, we cannot compute $\mathcal{R}(f)$ directly. Instead, we use the empirical risk as a surrogate,

$$\hat{\mathcal{R}}(f) \doteq \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)).$$

The goal is to obtain the empirical risk minimizer,

$$f^* \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathcal{R}}(f),$$

where \mathcal{F} is a family of functions that we assume f belongs to.

2 Anomaly detection

In anomaly detection, we are given a sample of objects $\mathcal{X} \subseteq \mathbb{R}^d$ with a normal class $\mathcal{N} \subseteq \mathcal{X}$ —the data points that we wish to keep. We wish to construct an anomaly classifier $\phi : \mathcal{X} \rightarrow \{0, 1\}$, such that

$$\phi(x) = 1 \iff x \notin \mathcal{N}.$$

Intuitively, an anomaly is an unlikely event. Hence, the strategy is to fit a model of a parametric family of distributions to the data \mathcal{X} ,

$$\mathcal{H} = \{p(\cdot \mid \theta) \mid \theta \in \Theta\}.$$

Then, we define the anomaly score of x as its probability $p(x \mid \theta^*)$ according to the optimal model in this hypothesis class. If x has a low score, it is more likely to be an anomaly.

Anomaly detection in a high-dimensional space is hard, because the normal class can be very complex. The idea is to project \mathcal{X} down to a lower dimensionality and perform anomaly detection there—hopefully the projected version of the normal class $\Pi(\mathcal{N})$ is less complex. In order to find the optimal linear projection, we will use PCA (*Principal Component Analysis*). Furthermore, it has been observed that linear projections of high-dimensional distributions onto low-dimensional spaces resemble Gaussian distributions. Hence, after performing PCA, we will fit a GMM (*Gaussian Mixture Model*).

2.1 Principal component analysis

The goal of PCA is to linearly project \mathbb{R}^d to \mathbb{R}^{d^-} such that the maximum amount of variance of the data is preserved.¹ Consider the base case $d^- = 1$. Let $u \in \mathbb{R}^d$ with $\|u\| = 1$, we project onto u by inner product,

$$x \mapsto u^\top x.$$

The sample mean of the reduced dataset is computed by

$$\frac{1}{n} \sum_{i=1}^n u^\top x_i = u^\top \bar{x},$$

where \bar{x} is the sample mean of the original dataset. Further, the sample variance of the reduced dataset is

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (u^\top x_i - u^\top \bar{x})^2 &= \frac{1}{n} \sum_{i=1}^n (u^\top (x_i - \bar{x}))^2 \\ &= \frac{1}{n} \sum_{i=1}^n u^\top (x_i - \bar{x})(x_i - \bar{x})^\top u \\ &= u^\top \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top \right) u \\ &= u^\top \Sigma u, \end{aligned}$$

¹ Components with larger variance are more informative.

where Σ is the covariance matrix of the dataset. Since we want the projection that preserves the maximum variance, we have the following objective,

$$\mathbf{u}^* \in \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u}.$$

The Lagrangian of this problem is

$$\mathcal{L}(\mathbf{u}; \lambda) = \mathbf{u}^\top \Sigma \mathbf{u} + \lambda(1 - \|\mathbf{u}\|^2)$$

with gradient

$$\frac{\partial \mathcal{L}(\mathbf{u}; \lambda)}{\partial \mathbf{u}} = 2\Sigma \mathbf{u} - 2\lambda \mathbf{u} \stackrel{!}{=} 0.$$

So, \mathbf{u} must satisfy $\Sigma \mathbf{u} = \lambda \mathbf{u}$ — \mathbf{u} is an eigenvector of Σ . It is easy to see that this must be the principal eigenvector by rewriting the objective,

$$\begin{aligned} \mathbf{u}^* &\in \operatorname{argmax}_{\|\mathbf{u}\|=1} \mathbf{u}^\top \Sigma \mathbf{u} \\ &= \operatorname{argmax}_{\substack{\|\mathbf{u}\|=1 \\ (\mathbf{u}, \lambda) \in \operatorname{eig}(\Sigma)}} \lambda \|\mathbf{u}\|^2 \\ &= \operatorname{argmax}_{\substack{\mathbf{u} \in \mathbb{R}^d \\ (\mathbf{u}, \lambda) \in \operatorname{eig}(\Sigma)}} \lambda \\ &= \mathbf{u}_1. \end{aligned}$$

For $d^- > 1$, the remaining principal components can be computed with a similar idea. Iteratively, we factor out the previous principal components and do as above on the transformed dataset. For example, to get the second principal component, we first factor out the first principal component,

$$\mathcal{X}_1 \doteq \{x - \operatorname{proj}_{\mathbf{u}_1}(x) \mid x \in \mathcal{X}\} = \{x - \mathbf{u}_1^\top x \cdot \mathbf{u}_1 \mid x \in \mathcal{X}\}.$$

Then, we do the same as above. While the above is easier to justify, we can also just take the top- d^- principal eigenvectors of $\Sigma = \frac{1}{n} \sum_{i=1}^n (x - \bar{x})(x - \bar{x})^\top$.

2.2 Gaussian mixture model

The PDF (*Probability Density Function*) of a GMM with k components is formalized as a convex combination of Gaussians,

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

For completeness, the multivariate Gaussian is defined as

$$\mathcal{N}(x; \mu, \Sigma) \doteq \frac{1}{(2\pi)^{n/2} \det(\Sigma)^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right).$$

So, the parameters of this model are

$$\theta = \{\pi_j, \mu_j, \Sigma_j \mid j \in [k]\},$$

```

1: Initialize  $\theta_0$ 
2: for  $t \in [T]$  do
3:    $q^* \in \operatorname{argmin}_q E(q, \theta_{t-1})$ 
4:    $\theta_t \in \operatorname{argmax}_\theta M(q^*, \theta)$ 
5: end for
6: return  $\theta_T$ 

```

where $\pi \in \Delta^{k-1}$ and $\{\Sigma_j \mid j \in [k]\}$ are positive definite. We fit the parameters of this model by maximizing the log-likelihood,

$$\begin{aligned}
\theta^* &\in \operatorname{argmax}_{\theta \in \Theta} \log p(\{x_1, \dots, x_n\}; \theta) \\
&= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log p(x_i; \theta) \\
&= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log \sum_{j=1}^k \pi_j \mathcal{N}(x_i; \mu_j, \Sigma_j).
\end{aligned}$$

Note that the above is intractable, so we would like to decompose it into tractable terms that can be computed. Let's assume that we know from which latent variable each data point was generated,² then we can compute the MLE of the extended dataset $\{(x_i, z_i) \mid i \in [n]\}$ as

$$\begin{aligned}
\log p(X, z; \theta) &= \sum_{i=1}^n \log p(x_i, z_i) \\
&= \sum_{i=1}^n \log(\pi_{z_i} \mathcal{N}(x_i; \mu_{z_i}, \Sigma_{z_i})) \\
&= \sum_{i=1}^n \log \pi_{z_i} + \sum_{i=1}^n \log \mathcal{N}(x_i; \mu_{z_i}, \Sigma_{z_i}),
\end{aligned}$$

which is tractable to maximize. Let q be a distribution over $[k]$, then we can rewrite the log-likelihood into tractable terms,

$$\begin{aligned}
\log p(X; \theta) &= \mathbb{E}_{z \sim q} [\log p(X; \theta)] \\
&= \mathbb{E}_{z \sim q} \left[\log \left(\frac{p(X, z; \theta)}{p(z \mid X; \theta)} \right) \right] \\
&= \mathbb{E}_{z \sim q} \left[\log \left(\frac{p(X, z; \theta)}{p(z \mid X; \theta)} \frac{q(z)}{q(z)} \right) \right] \\
&= \underbrace{\mathbb{E}_{z \sim q} \left[\log \frac{p(X, z; \theta)}{q(z)} \right]}_{\doteq M(q, \theta)} + \underbrace{\mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z \mid X; \theta)} \right]}_{\doteq E(q, \theta)}.
\end{aligned}$$

These terms have the following two properties,

$$\begin{aligned}
\log p(X; \theta) &\geq M(q, \theta), \quad \forall q, \theta \\
\log p(X; \theta) &= M(q^*, \theta), \quad q^* = p(\cdot \mid X; \theta), \quad \forall \theta.
\end{aligned}$$

Hence, we can use $M(q^*, \theta)$ as an approximation of $\log p(X; \theta)$ if $q = q^*$. So, we alternately optimize $M(q, \theta)$ w.r.t. θ and $E(q, \theta)$ w.r.t. q . In the

Algorithm 1. The expectation-maximization algorithm, where

$$\begin{aligned}
M(q, \theta) &\doteq \mathbb{E}_{z \sim q} \left[\log \frac{p(X, z; \theta)}{q(z)} \right] \\
E(q, \theta) &\doteq \mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p(z \mid X; \theta)} \right].
\end{aligned}$$

² The Gaussians are the latent variables here. The generative story goes that the data is sampled by first sampling a Gaussian from π and then the Gaussian is sampled. Here, we assume that we know from which Gaussian each point was sampled.

$$E(q, \theta) \geq 0.$$

In this case, it is easy to see that $E(q^*, \theta) = 0$.

case of a Gaussian mixture model, we can find the optima in a closed form.

Theorem 2.1 (EM algorithm convergence). Using the expectation-maximization algorithm, $\{\log p(x; \theta_t)\}_{t=0}^T$ is non-decreasing.

Proof. Given a data point x and current parameters θ , we have the following update,

$$\theta' \in \operatorname{argmax}_{\theta \in \Theta} M(q^*, \theta).$$

Hence, we have

$$\log p(x; \theta) = M(q^*, \theta) \leq M(q^*, \theta') \leq \log p(x; \theta').$$

Thus, $\{\log p(x; \theta_t)\}_{t=0}^T$ is non-decreasing. ■

2.3 Summary

In conclusion, given a set of data points \mathcal{X} with normal points $\mathcal{N} \subseteq \mathcal{X}$, we train an anomaly detector as follows,

1. Fit a projector $\pi : \mathbb{R}^d \rightarrow \mathbb{R}^{d-}$ using PCA;
2. Fit a GMM with k components $p(\cdot | \theta)$ to $\{\pi(x) | x \in \mathcal{X}\}$ using the EM algorithm;
3. For a point $x \in \mathcal{X}$, its “anomaly score” is computed by $-\log p(\pi(x); \theta)$.

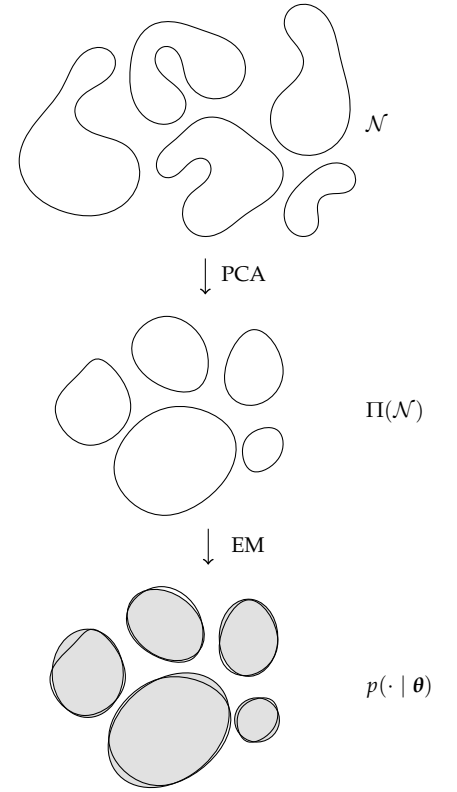


Figure 2.1. Illustration of the anomaly detection pipeline.

3 Density estimation

In this section, we will consider parametric models,

$$\{p(\cdot; \theta) \mid \theta \in \Theta\}.$$

The problem we concern ourselves with is finding the best parameter θ .

The most common method of finding the best parameters is MLE,

$$\begin{aligned}\hat{\theta}_{\text{MLE}} &\in \operatorname{argmax}_{\theta \in \Theta} \prod_{i=1}^n p(x_i; \theta) \\ &= \operatorname{argmin}_{\theta \in \Theta} - \sum_{i=1}^n \log p(x_i; \theta).\end{aligned}$$

The following properties makes the MLE estimator attractive,

1. *Consistency*: In the limit of n , $\hat{\theta}_{\text{MLE}}$ converges to the true parameter θ^* ;
2. *Equivariance*: If $\hat{\theta}$ is the MLE of θ , then $g(\hat{\theta})$ is the MLE of $g(\theta)$ for any function g ;
3. *Asymptotically normal*: In the limit of n , $\hat{\theta}_{\text{MLE}} - \theta / \sqrt{n}$ converges to a random variable with distribution $\mathcal{N}(0, \mathcal{I}(\theta)^{-1})$, where \mathcal{I} is the Fisher information matrix;
4. *Asymptotically efficient*: In the limit of n , the MLE estimator has the smallest variance of all unbiased estimators.

We can understand the asymptotic efficiency of estimators better by the Rao–Cramér bound, which provides a bound on the variance of the estimator. We will only consider the general scalar case, but it generalizes to the multivariate case.

Theorem 3.1 (Rao–Cramér bound (scalar case)). For any (scalar) unbiased estimator $\hat{\theta} : \mathcal{Y}^n \rightarrow \mathbb{R}$, given n data points, of $\theta \in \mathbb{R}$, we have the following bound on its variance,

$$\operatorname{Var}[\hat{\theta}(\mathbf{y})] \geq \frac{\left(\frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1\right)^2}{\mathcal{I}_n(\theta)} + b_{\hat{\theta}}^2,$$

where $\mathcal{I}_n(\theta)$ is the Fisher information,

$$\mathcal{I}_n(\theta) \doteq \mathbb{E}_{\mathbf{y}|\theta} \left[\left(\frac{\partial}{\partial \theta} \log p(\mathbf{y} \mid \theta) \right)^2 \right] \stackrel{i.i.d.}{=} n \cdot \mathcal{I}(\theta),$$

and $b_{\hat{\theta}}$ is the bias of $\hat{\theta}$,

$$b_{\hat{\theta}} \doteq \mathbb{E}_{\mathbf{y}|\theta} [\hat{\theta}(\mathbf{y})] - \theta.$$

Proof. Let the “score” be defined as follows,

$$\Lambda(\mathbf{y}, \theta) \doteq \frac{\partial}{\partial \theta} \log p(\mathbf{y} \mid \theta) = \frac{1}{p(\mathbf{y} \mid \theta)} \frac{\partial}{\partial \theta} p(\mathbf{y} \mid \theta).$$

Then we have $\mathcal{I}_n(\theta) = \mathbb{E}_{\mathbf{y}|\theta} [\Lambda(\mathbf{y}, \theta)^2]$.

The expected score is equal to zero,

$$\begin{aligned}
 \mathbb{E}_{\mathbf{y}|\theta}[\Lambda(\mathbf{y}, \theta)] &= \int p(\mathbf{y} | \theta) \Lambda(\mathbf{y}, \theta) d\mathbf{y} \\
 &= \int \frac{\partial}{\partial \theta} p(\mathbf{y} | \theta) d\mathbf{y} \\
 &= \frac{\partial}{\partial \theta} \int p(\mathbf{y} | \theta) d\mathbf{y} \\
 &= \frac{\partial}{\partial \theta} 1 \\
 &= 0.
 \end{aligned}$$

(Hence, the Fisher information is equivalent to the variance of the score.)

Furthermore, the cross-correlation between the score $\Lambda(\mathbf{y}, \theta)$ and the estimator $\hat{\theta}(\mathbf{y})$ can be computed as

$$\begin{aligned}
 \text{Cov}_{\mathbf{y}|\theta}(\Lambda(\mathbf{y}, \theta), \hat{\theta}(\mathbf{y})) &= \mathbb{E}_{\mathbf{y}|\theta}[(\Lambda(\mathbf{y}, \theta) - \mathbb{E}[\Lambda(\mathbf{y}, \theta)]) (\hat{\theta}(\mathbf{y}) - \mathbb{E}[\hat{\theta}(\mathbf{y})])] \\
 &= \mathbb{E}_{\mathbf{y}|\theta}[\Lambda(\mathbf{y}, \theta) \hat{\theta}(\mathbf{y})] - \mathbb{E}_{\mathbf{y}|\theta}[\Lambda(\mathbf{y}, \theta)] \mathbb{E}_{\mathbf{y}|\theta}[\hat{\theta}(\mathbf{y})] \\
 &= \mathbb{E}_{\mathbf{y}|\theta}[\Lambda(\mathbf{y}, \theta) \hat{\theta}(\mathbf{y})] \\
 &= \int p(\mathbf{y} | \theta) \Lambda(\mathbf{y}, \theta) \hat{\theta}(\mathbf{y}) d\mathbf{y} \\
 &= \int \frac{\partial}{\partial \theta} p(\mathbf{y} | \theta) \hat{\theta}(\mathbf{y}) d\mathbf{y} \\
 &= \frac{\partial}{\partial \theta} \int p(\mathbf{y} | \theta) \hat{\theta}(\mathbf{y}) d\mathbf{y} \\
 &= \frac{\partial}{\partial \theta} (\mathbb{E}_{\mathbf{y}|\theta}[\hat{\theta}(\mathbf{y})] - \theta) + 1 \\
 &= \frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1.
 \end{aligned}$$

Using the Cauchy-Schwarz inequality, we have the following bound,

$$\begin{aligned}
 \left(\frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1 \right)^2 &= \text{Cov}_{\mathbf{y}|\theta}(\Lambda(\mathbf{y}, \theta), \hat{\theta}(\mathbf{y}))^2 \\
 &= \left(\mathbb{E}_{\mathbf{y}|\theta}[\Lambda(\mathbf{y}, \theta) (\hat{\theta}(\mathbf{y}) - \mathbb{E}[\hat{\theta}(\mathbf{y})])] \right)^2 \\
 &\leq \mathbb{E}_{\mathbf{y}|\theta}[\Lambda(\mathbf{y}, \theta)^2] \cdot \mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \mathbb{E}[\hat{\theta}(\mathbf{y})])^2] \\
 &= \mathcal{I}_n(\theta) \cdot \mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \theta - \mathbb{E}[\hat{\theta}(\mathbf{y}) - \theta])^2] \\
 &= \mathcal{I}_n(\theta) \cdot \left(\mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \theta)^2] + (\mathbb{E}_{\mathbf{y}|\theta}[\hat{\theta}(\mathbf{y})] - \theta)^2 \right. \\
 &\quad \left. - 2\mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \theta)]^2 \right) \\
 &= \mathcal{I}_n(\theta) \cdot \left(\mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \theta)^2] - (\mathbb{E}_{\mathbf{y}|\theta}[\hat{\theta}(\mathbf{y})] - \theta)^2 \right) \\
 &= \mathcal{I}_n(\theta) \cdot \left(\mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \theta)^2] - b_{\hat{\theta}}^2 \right).
 \end{aligned}$$

Cauchy-Schwarz inequality.

Re-arranging yields

$$\text{Var}_{\mathbf{y}|\theta}[\hat{\theta}(\mathbf{y})] \doteq \mathbb{E}_{\mathbf{y}|\theta}[(\hat{\theta}(\mathbf{y}) - \theta)^2] \geq \frac{\left(\frac{\partial}{\partial \theta} b_{\hat{\theta}} + 1 \right)^2}{\mathcal{I}_n(\theta)} + b_{\hat{\theta}}^2.$$



Note the trade-off between variance and bias. If $\frac{\partial}{\partial \theta} b_{\hat{\theta}} < 0$, then the variance might be smaller than the variance of the best unbiased estimator.

Corollary. Let $\hat{\theta}$ be an unbiased estimator—i.e., $b_{\hat{\theta}} = 0$ —then

$$\text{Var}[\hat{\theta}(y)] \geq \frac{1}{\mathcal{I}_n(\theta)}.$$

Lemma 3.2. The maximum likelihood estimator $\hat{\theta}_{\text{ML}}$ is asymptotically efficient,

$$\lim_{n \rightarrow \infty} \text{Var}[\hat{\theta}_{\text{ML}}] = \frac{1}{\mathcal{I}_n(\theta)}.$$

However, the MLE estimator is not necessarily efficient for finite samples.

4 Regression

In regression, we try to estimate a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that best fits a given dataset $\{(x_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^d \times \mathbb{R}$. I.e., we try to minimize the following loss function,

$$\ell(f) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2.$$

4.1 Linear regression

In linear regression, we assume the underlying data to be linear with a fixed noise,

$$Y | X = x \sim \mathcal{N}(\beta_*^\top x, \sigma^2)$$

for some ground truth β_* . We parametrize f as a linear function,

$$f(x; \beta) = \beta^\top x.$$

Under these assumptions, the minimizer of the loss function is analytically tractable—the OLSE (*Ordinary Least Squares Estimator*),

$$\hat{\beta} = (X^\top X)^{-1} X^\top y,$$

where $X \in \mathbb{R}^{n \times d}$ is the design matrix with respective outputs $y \in \mathbb{R}^n$.

In practice, it is important to remove outliers, because linear models can be heavily influenced by them. Also, the data should be standardized, such that all features are on the same scale, because differences in scale can make matrix inversion unstable.

Moreover, it has been found that in high-dimensional problems, logistic regression produces overconfident outputs [Sur and Candès, 2019]. For example, see the 20-dimensional β^* in Figure 4.1. If we sample many datapoints with some noise and train a logistic regression model on them, the colored $\hat{\beta}$ in Figure 4.1 is produced. As can be seen, the weights of the first 10 dimensions are overestimated. Because of this, the model becomes overconfident as shown in Figure 4.2. This is called “the curse of dimensionality”.

Further, if features are collinear, then the model might only learn the correlation with the target variable of one of them and discard the other. In addition, some singular values will equal zero. This makes matrix inversion unstable. This can easily be shown by considering the singular

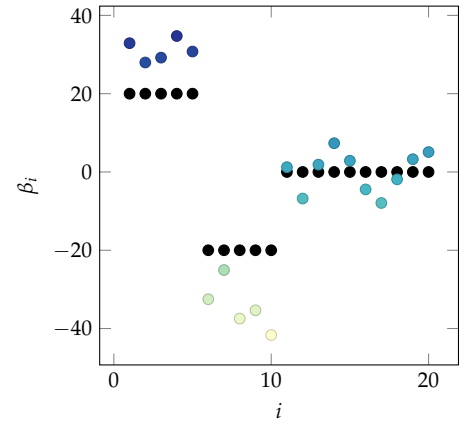


Figure 4.1. β^* is shown as black marks and $\hat{\beta}$ is indicated by the marks. As can be seen, $\hat{\beta}$ is overestimated for indices where $\beta_i^* \neq 0$.

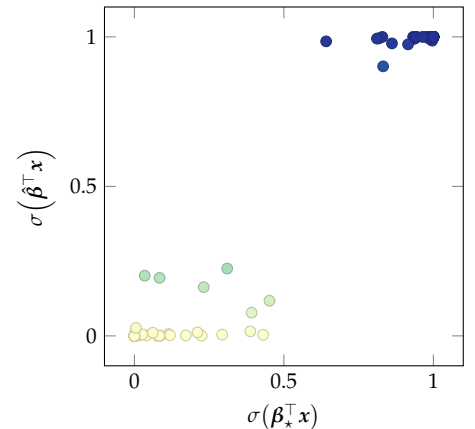


Figure 4.2. True vs. predicted probability—the estimator is overconfident in its predictions.

value decomposition $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$, then

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{V}\mathbf{D}^\top \mathbf{U}^\top \mathbf{U}\mathbf{D}\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{D}^\top \mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}\mathbf{D}^\top \mathbf{D}\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{D}^\top \mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}(\mathbf{D}^{-1})^\top \mathbf{D}^{-1} \mathbf{V}^\top \mathbf{V}\mathbf{D}^\top \mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}(\mathbf{D}^{-1})^\top \mathbf{D}^{-1} \mathbf{D}^\top \mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}\mathbf{D}^{-1} \mathbf{U}^\top \mathbf{y}.\end{aligned}$$

The inversion of \mathbf{D} is unstable if the singular values are small. The solution to this is to remove collinear features and/or to add regularization.

Lemma 4.1 (Risk decomposition). Consider the mean-squared error loss. The risk (expected loss) can be decomposed as follows into bias, variance, and noise terms,

$$\mathbb{E} \left[(\hat{f}(\mathbf{x}) - y)^2 \right] = \underbrace{(\mathbb{E}[\hat{f}(\mathbf{x})] - \mathbb{E}[y])^2}_{\text{squared bias}} + \underbrace{\text{Var}[\hat{f}(\mathbf{x})]}_{\text{variance}} + \underbrace{\text{Var}[y]}_{\text{noise}}.$$

Proof.

$$\begin{aligned}\mathbb{E} \left[(\hat{f}(\mathbf{x}) - y)^2 \right] &= \mathbb{E} \left[(\hat{f}(\mathbf{x}) - f(\mathbf{x}) + \epsilon)^2 \right] \\ &= \mathbb{E} \left[(\hat{f}(\mathbf{x}) - f(\mathbf{x}))^2 \right] + \mathbb{E}[\epsilon^2] + 2\mathbb{E}[\epsilon]\mathbb{E}[\hat{f}(\mathbf{x}) - f(\mathbf{x})] \\ &= \mathbb{E} \left[(\hat{f}(\mathbf{x}) - f(\mathbf{x}))^2 \right] + \text{Var}[y] \\ &= \mathbb{E} \left[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})] + \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}))^2 \right] + \text{Var}[y] \\ &= \mathbb{E} \left[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2 \right] + \mathbb{E}[\mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x})]^2 \\ &\quad + 2\mathbb{E} \left[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})]) (f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})]) \right] + \text{Var}[y] \\ &= (\mathbb{E}[\hat{f}(\mathbf{x}) - \mathbb{E}[y]])^2 + \text{Var}[\hat{f}(\mathbf{x})] + \text{Var}[y].\end{aligned}$$

$$y = f(\mathbf{x}) + \epsilon, \text{ where } \mathbb{E}[\epsilon] = 0.$$

$$\mathbb{E}[\epsilon^2] = \mathbb{E}[(y - f(\mathbf{x}))^2] = \mathbb{E}[(y - \mathbb{E}[y])^2] = \text{Var}[y].$$

$f(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})]$ is deterministic so we can pull it out, then the term becomes 0. Also, $f(\mathbf{x}) = \mathbb{E}[y]$.

Theorem 4.2 (Gauss-Markov). The ordinary least-squares estimator is the (unique) minimum-variance unbiased linear estimator.

Just because the OLSE is the minimum variance unbiased estimator, it does not mean it is the best, because introducing some bias may considerably decrease the variance. Bayesianism adds bias by introducing a prior, which often induces a regularizing term that decreases variance.

4.2 Regularization

By introducing a Gaussian prior, $\beta \sim \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$, the MAP estimator is the Ridge regression estimator. It essentially just adds an ℓ_2 -norm regularization with weight λ . Furthermore, we could also use a Laplacian prior, $\beta \sim \text{Lap}(\mathbf{0}, \lambda \mathbf{I})$, which results in the LASSO estimator that adds an ℓ_1 -norm regularization with weight λ . The ℓ_1 -norm results in sparse features which can be interpreted better. For example, it has the feature that as λ increases, weights do not go from positive to negative or vice versa. This is not the case for the ℓ_2 -norm.

4.3 Polynomial regression

In polynomial regression, we construct a feature function of all possible polynomials, e.g.,

$$\phi([x_1, x_2]) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2, x_1^3, x_1^2 x_2, \dots].$$

We then perform linear regression in this space,

$$\psi(x; \beta) = \langle \beta, \phi(x) \rangle.$$

The problem is that this space is infinitely dimensional, and the inner product is ill-defined in an infinitely dimensional space.³ We can solve this by transforming the vector such that inner products cannot diverge by introducing a data-dependent scalar,

$$\phi(x) = \alpha(x) \cdot \tilde{\phi}(x), \quad \alpha(x) > 0.$$

³ This is due to $\langle \phi(x), \phi(x') \rangle = \infty$ for some x, x' , e.g., $x = x' = [1, 1]$.

$\tilde{\phi}$ contains all polynomials in this equation.

Specifically,

$$\phi(x) = \exp\left(-\frac{1}{2}\|x\|^2\right) \left[\frac{\prod_{i=1}^d x_i^{\alpha_i}}{\sqrt{\prod_{i=1}^d \alpha_i!}} \right]_{\alpha \in \mathbb{N}^d}.$$

An inner product w.r.t. this feature function is a valid inner product. It results in the RBF (*Radial Basis Function*) kernel,

$$\begin{aligned} \langle \phi(x), \phi(x') \rangle &= \exp\left(-\frac{1}{2}\|x\|^2\right) \exp\left(-\frac{1}{2}\|x'\|^2\right) \sum_{\alpha \in \mathbb{N}^d} \frac{\prod_{i=1}^d x_i^{\alpha_i} (x'_i)^{\alpha_i}}{\prod_{i=1}^d \alpha_i!} \\ &= \exp\left(-\frac{1}{2}\|x\|^2\right) \exp\left(-\frac{1}{2}\|x'\|^2\right) \exp(x^\top x') \\ &= \exp\left(-\frac{1}{2}(\|x\|^2 + \|x'\|^2 - 2\langle x, x' \rangle)\right) \\ &= \exp\left(-\frac{1}{2}\|x - x'\|^2\right). \end{aligned}$$

Using the multinomial series expansion.

Cosine theorem.

Let $\Phi \in \mathbb{R}^{n \times \infty}$ contain all feature vectors, then the OLSE for polynomial regression can be computed by

$$\hat{\beta} = (\Phi^\top \Phi)^{-1} \Phi^\top y, \quad \Phi \in \mathbb{R}^{n \times \infty}.$$

However, $\Phi^\top \Phi$ cannot be computed, because it is $\infty \times \infty$ -dimensional. We can fix this by observing that $\Phi \Phi^\top \in \mathbb{R}^{n \times n}$ and rewriting the OLSE as

$$\begin{aligned}\hat{\beta} &= (\Phi^\top \Phi)^{-1} \Phi^\top y \\ &= (\Phi^\top \Phi)^{-1} \Phi^\top (\Phi \Phi^\top) (\Phi \Phi^\top)^{-1} y \\ &= (\Phi^\top \Phi)^{-1} \Phi^\top \Phi \Phi^\top (\Phi \Phi^\top)^{-1} y \\ &= \Phi^\top (\Phi \Phi^\top)^{-1} y.\end{aligned}$$

Let $K = \Phi \Phi^\top$, then it only contains kernel evaluations— $k_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$. The next problem is that Φ^\top is still infinitely dimensional. However, this is not a problem, since when we want to make a prediction, we do the following,

$$\psi(\hat{x}) = \langle \phi(\hat{x}), \hat{\beta} \rangle = \phi(\hat{x})^\top \Phi^\top K^{-1} y.$$

Let $k(\hat{x}) = \phi(\hat{x})^\top \Phi^\top$, then it only contains kernel evaluations with the new point— $k_i(\hat{x}) = \langle \phi(\hat{x}), \phi(x_i) \rangle$. In conclusion, we can make predictions by

$$\psi(\hat{x}) = k(\hat{x})^\top K^{-1} y, \quad k(\hat{x}) \in \mathbb{R}^n, K \in \mathbb{R}^{n \times n}, y \in \mathbb{R}^n.$$

However, the problem with this approach is that it takes $O(n^3)$ to make a prediction, because it scales in the number of data points. Furthermore, we need to store all data points—we cannot compress them into model parameters.

5 Causality

In general, models do not distinguish between causal and non-causal factors in the feature space. Therefore, they might identify non-causal factors as highly correlating with the output variable. *E.g.*, when classifying images as either depicting a cow or a camel, the model might identify the background as an important feature, because cows are usually in grass and camels are usually in the desert. However, if this classifier were to see a cow with a different background, it would fail.

The following are examples of causal fallacies, where one might conclude that X causes Y ,

- Reverse causality, where Y actually predicts X and not the other way around;
- Third-cause fallacy, where there is a confounding factor Z that influences both X and Y ;
- Bidirectional causation, where X and Y influence each other.

Further, a domain shift happens when the test samples are drawn from a different distribution than the training samples. *E.g.*, COVID-19 detection models trained on a western population might not perform well on an eastern population.

Moreover, shortcut learning happens when there is a spurious correlation between causal and non-causal features in the training dataset. The resulting estimator abuses the non-causal features to generalize in the testing dataset. The solution to this is to encode the features such that they do not depend on the environment.

5.1 Counterfactual invariance

Let X be the feature vector representing the object and let Y be a target variable of interest, and let f be the function that estimates Y from X . Further, let W describe features that influence X , but should not influence the estimator f . Let a counterfactual be denoted as $X(w)$, which is the feature vector we would have obtained if we would have had $W = w$, after the fact. Then, the estimator f is counterfactually invariant if

$$f(X(w)) = f(X(w')), \quad \forall w, w' \in \text{range}(W).$$

In words, f is not influenced by the value of W . One way of obtaining a counterfactually invariant estimator is by extending the training dataset to contain enough counterfactuals. This can, for example, be achieved by data augmentation. However, this is not always possible.

Let X_A be the parts of X causally influenced by A , and let A^\perp be the set of variables independent of A . Then, f is counterfactually invariant if and only if f only depends on X_{W^\perp} .

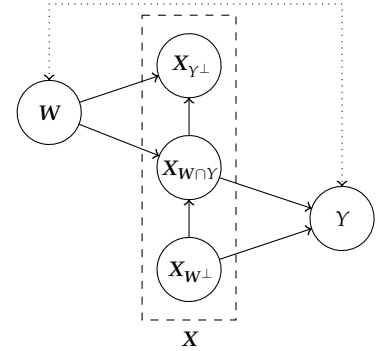


Figure 5.1. Causal graph of W , X , and Y . X is decomposed into several components. Dotted arrows are non-causal relationships.

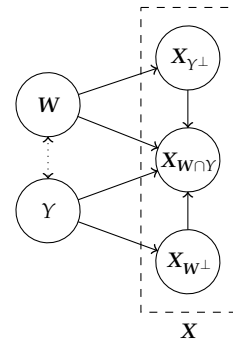


Figure 5.2. Anti-causal graph of W , X , and Y . X is decomposed into several components. Dotted arrows are non-causal relationships.

Theorem 5.1 (Necessary conditions for counterfactual invariance [Veitch et al., 2021]). If f is a counterfactually invariant predictor, then

- In the anti-causal scenario, $f(\mathbf{X}) \perp \mathbf{W} \mid Y$;
- In the causal scenario without selection (but possibly confounded), $f(\mathbf{X}) \perp \mathbf{W}$;
- In the causal scenario without confoundedness (but possibly with selection), $f(\mathbf{X}) \perp \mathbf{W} \mid Y$, as long as $\mathbf{X} \perp Y \mid \mathbf{X}_{\mathbf{W}^\perp}, \mathbf{W}$.

Proof. This can be easily be seen by using d-separation—see Figures 5.1 to 5.3. ■

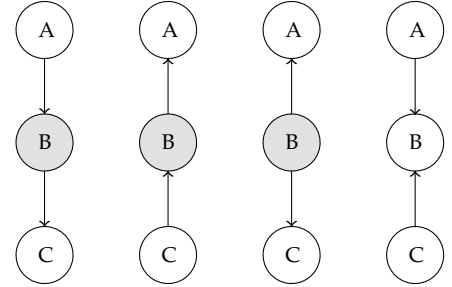


Figure 5.3. d-separation. Here, we condition (or observe) on nodes that are gray. In all of these examples A is separated from C , meaning that $A \perp C \mid Z$, where Z are the conditioned variables.

6 Gaussian processes

Let the inputs be $\mathbf{X} \in \mathbb{R}^{n \times d}$, the outputs be $\mathbf{y} \in \mathbb{R}^n$, and the weight matrix be $\boldsymbol{\beta} \in \mathbb{R}^d$, then linear regression models the generative process as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

This is equivalent to defining a Gaussian over the predictions,

$$\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}).$$

BLR (*Bayesian Linear Regression*) extends linear regression by defining a prior over the regression coefficients,

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}^{-1}), \quad \boldsymbol{\Lambda} \in \mathbb{R}^{d \times d}.$$

The posterior can be analytically computed as

$$\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where

$$\boldsymbol{\mu} = \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}, \quad \boldsymbol{\Sigma} = \sigma^2 (\mathbf{X}^\top \mathbf{X} + \sigma^2 \boldsymbol{\Lambda})^{-1}.$$

Proof.

$$\begin{aligned} p(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}) &\propto p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}) \cdot p(\boldsymbol{\beta}) && \text{Bayes' rule.} \\ &= \mathcal{N}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}) \cdot \mathcal{N}(\boldsymbol{\beta}; \mathbf{0}, \boldsymbol{\Lambda}^{-1}) \\ &\propto \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} (\|\mathbf{y}\|^2 + \|\mathbf{X}\boldsymbol{\beta}\|^2 - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta}) + \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta} \right)\right) && \text{Cosine theorem.} \\ &= \exp\left(-\frac{1}{2} \left(\frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{y} + \frac{1}{\sigma^2} \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} - \frac{2}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta}^\top \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \boldsymbol{\Lambda} \right) \boldsymbol{\beta} - \frac{2}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} + \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{y} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} - \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} + \frac{1}{\sigma^2} \mathbf{y}^\top \mathbf{y} \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta}^\top \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right) - \frac{1}{\sigma^2} (\boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y})^\top \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right) \right)\right) \\ &= \exp\left(-\frac{1}{2} \left(\left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right)^\top \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{\beta} - \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y} \right) \right)\right) \\ &\propto \mathcal{N}\left(\boldsymbol{\beta}; \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{y}, \boldsymbol{\Sigma}\right). \end{aligned}$$

■

In general, we do not know the true weights that generated the data points. But, we can now define a joint distribution over output variables with unknown weights,

$$\mathbf{y} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{X} \boldsymbol{\Lambda}^{-1} \mathbf{X}^\top + \sigma^2 \mathbf{I}).$$

Proof. As we saw earlier, we compute outputs as follows,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

However, now the weights are unknown, so we make use of its prior to compute the distribution over \mathbf{y} ,

$$\begin{aligned} \mathbb{E}[\mathbf{y}] &= \mathbb{E}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}] \\ &= \mathbf{X}\mathbb{E}[\boldsymbol{\beta}] + \mathbb{E}[\boldsymbol{\epsilon}] \\ &= \mathbf{0}. \end{aligned}$$

$$\begin{aligned} \text{Cov}[\mathbf{y}] &= \mathbb{E}[(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})^\top] \\ &= \mathbb{E}[\mathbf{X}\boldsymbol{\beta}\boldsymbol{\beta}^\top \mathbf{X}^\top] + \mathbb{E}[\mathbf{X}\boldsymbol{\beta}\boldsymbol{\epsilon}^\top] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\beta}^\top \mathbf{X}^\top] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] \\ &= \mathbf{X}\mathbb{E}[\boldsymbol{\beta}\boldsymbol{\beta}^\top] \mathbf{X}^\top + \mathbf{X}\mathbb{E}[\boldsymbol{\beta}]\mathbb{E}[\boldsymbol{\epsilon}^\top] + \mathbb{E}[\boldsymbol{\epsilon}]\mathbb{E}[\boldsymbol{\beta}^\top] \mathbf{X}^\top + \sigma^2 \mathbf{I} \\ &= \mathbf{X}\boldsymbol{\Lambda}^{-1} \mathbf{X}^\top + \sigma^2 \mathbf{I}. \end{aligned}$$

$\boldsymbol{\beta}$ and $\boldsymbol{\epsilon}$ are independent.

■

GPs (*Gaussian Processes*) generalize BLR by observing that we can kernelize the covariance matrix,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \boldsymbol{\Lambda}^{-1} \mathbf{x}_j.$$

We could instead use any other kernel function to model other functions.

6.1 Kernels

Kernel functions specify the similarity between any two data points. They encode assumptions about the function that is to be learned.

Definition 6.1 (Kernel function). A kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ over $\mathcal{X} \subset \mathbb{R}^d$ must satisfy the following properties,

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= k(\mathbf{x}', \mathbf{x}) \\ \int \int k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' &\geq 0, \quad \forall f \in L_2. \end{aligned}$$

Symmetry.

Positive semi-definiteness.

Definition 6.2 (Stationary and isotropic). A kernel $k(\mathbf{x}, \mathbf{x}')$ is stationary if it only depends on $\mathbf{x} - \mathbf{x}'$. Further, it is isotropic if it only depends on $\|\mathbf{x} - \mathbf{x}'\|_2$.

The following are common kernels,

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$$

Linear kernel.

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + 1)^p, \quad p \in \mathbb{N}$$

Polynomial kernel.

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{\ell^2}\right), \quad \ell \in \mathbb{R}$$

RBF (*Radial Basis Function*) kernel.

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\kappa \mathbf{x}^\top \mathbf{x}') - b, \quad \kappa, b \in \mathbb{R}$$

Sigmoid kernel.

Different kernels have different invariance properties, such as invariance to rotation or translation. In order to learn invariances from data, many samples are needed. So, it might be better to encode them if we know them a priori.

Given two kernel functions k_1, k_2 defined on the same data space and positive scalar $c > 0$, the following are also valid kernels,

$$k(x, x') = k_1(x, x') + k_2(x, x')$$

$$k(x, x') = k_1(x, x') \cdot k_2(x, x')$$

$$k(x, x') = c \cdot k_1(x, x')$$

$$k(x, x') = \exp(k_1(x, x')).$$

In practice, the kernels are often composed together and hyperparameters are determined by performance on a held out validation dataset.

6.2 Prediction

As we saw earlier in the case of BLR, the marginal is jointly Gaussian,

$$\begin{bmatrix} y \\ y^* \end{bmatrix} \mid X, x^* \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K + \sigma^2 I & k(X, x^*) \\ k(x^*, X) & k(x^*, x^*) \end{bmatrix}\right).$$

Theorem 6.3 (Conditional Gaussian distribution). Given

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right).$$

Then, the conditional Gaussian is given by

$$x_2 \mid x_1 = z \sim \mathcal{N}\left(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(z - \mu_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right).$$

Using the above theorem, we can easily compute the predictive distribution,

$$y^* \mid x^*, X, y \sim \mathcal{N}\left(k^\top (K + \sigma^2 I)^{-1} y, c - k^\top (K + \sigma^2 I)^{-1} k\right),$$

where

$$K = K(X, X), \quad k = k(X, x^*), \quad c = k(x^*, x^*).$$

Now, we can compute a prediction along with its uncertainty. However, the problem with this approach is that it has $\mathcal{O}(n^3)$ runtime.

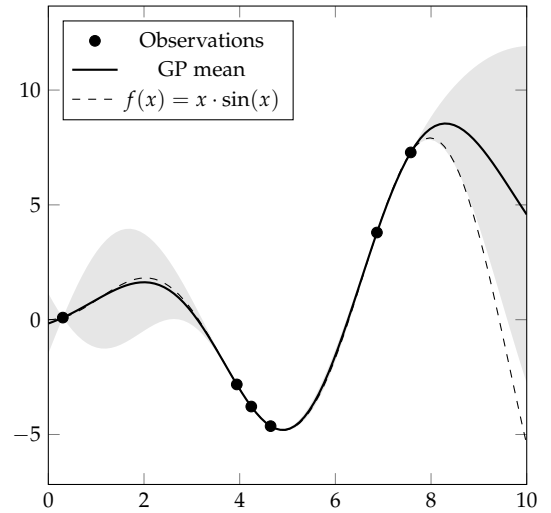


Figure 6.1. Fitted Gaussian process (RBF kernel, $\ell = 1.86$) with its 95% confidence interval.

7 Uncertainty quantification

7.1 Statistical model validation

In statistical modeling, we typically want to predict a target random variable $Y \in \mathcal{Y}$ given a random vector $\mathbf{X} \in \mathcal{X}$. Formally, we want to find a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ that minimizes the expected risk,

$$\mathcal{R}(f) \doteq \mathbb{E}_{\mathbf{X}, Y}[\mathbb{1}\{f(\mathbf{X}) \neq Y\}].$$

However, this function is intractable, because (1) we do not have access to the joint distribution over \mathbf{X}, Y , (2) it is intractable to find f without any assumptions on its structure, and (3) it is unclear how to minimize the expectation of the indicator function. To resolve (1), we obtain a dataset of samples from the joint distribution. Problem (2) can be resolved by restricting the functions to a parameterized hypothesis space \mathcal{H} . Lastly, the solution to (3) is to approximate the indicator function by a differentiable loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$. As a result, given a dataset $\mathcal{Z} = \{(x_i, y_i)\}_{i=1}^n$, we can approximate the expected risk by the empirical risk,

$$\hat{\mathcal{R}}(f) \doteq \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)).$$

We denote the minimizer of this function by \hat{f} .

Since \hat{f} is obtained by training on the dataset \mathcal{Z} , we need to find a way to measure the performance of \hat{f} on unseen data points—we want to estimate $\mathcal{R}(\hat{f})$. Furthermore, let $\mathcal{A} : \mathcal{P}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{H}$ be the algorithm that maps datasets to functions, then we want to estimate its expected risk,

$$\mathcal{R}(\mathcal{A}) \doteq \mathbb{E}_{\mathcal{Z}}[\mathcal{R}(\mathcal{A}(\mathcal{Z}))].$$

Cross-validation. Cross-validation works as follows,

1. Partition the data \mathcal{Z} in K equally sized disjoint subsets, such that

$$\mathcal{Z} = \bigcup_{k=1}^K \mathcal{Z}_k;$$

2. Use \mathcal{A} to produce K estimators \hat{f}^{-k} from $\mathcal{Z} - \mathcal{Z}_k$ for all $k \in [K]$;
3. Estimate the expected risk by

$$\mathcal{R}^{\text{CV}}(\mathcal{A}) \doteq \frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{f}^{-(k(i))}(x_i)),$$

where $k : [n] \rightarrow [K]$ maps i to the partition such that $x_i \in \mathcal{Z}_{k(i)}$.

In conclusion, the loss terms in the expected risk depend on approximators that have not seen the data it is evaluating.

Bootstrap. The bootstrap method is used for measuring the distribution over statistical parameters. It does so by creating B bootstrap samples by sampling with replacement from the dataset. For each bootstrap sample, the parameters are computed. Using these parameters, we can estimate the variability and uncertainty of the parameter. Note that this works for any type of parameter that can be estimated from the dataset—not only model parameters. Let $S : \mathcal{P}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}^d$ be a function that processes datasets and outputs parameters, then we can compute the statistics of this parameter as follows,

1. Draw B samples \mathcal{Z}^{*b} of size n from \mathcal{Z} with replacement;
2. Compute an estimate $S(\mathcal{Z}^{*b})$ for each bootstrap sample.

Now, we can use this set of parameter estimates to compute its statistics. For example, we can compute its mean and covariance,

$$\begin{aligned}\mu(S) &= \frac{1}{B} \sum_{b=1}^B S(\mathcal{Z}^{*b}) \\ \Sigma(S) &= \frac{1}{B-1} \sum_{b=1}^B \left(S(\mathcal{Z}^{*b}) - \mu(S) \right) \left(S(\mathcal{Z}^{*b}) - \mu(S) \right)^\top.\end{aligned}$$

The algorithm \mathcal{A} is a function that processes a dataset and outputs an estimator, which is used to compute a loss term. We can use this to estimate the empirical risk across n data points,

$$\hat{\mathcal{R}}^{\text{BS}}(\mathcal{A}) \doteq \frac{1}{n \cdot B} \sum_{b=1}^B \sum_{i=1}^n \ell(y_i, f^{*b}(x_i)).$$

However, computing this validation metric includes data that the bootstrap estimates were trained on—the estimate is overly optimistic. Thus, we can use the bootstrap method to estimate the expected risk $\mathcal{R}(\mathcal{A})$ as follows,

$$\mathcal{R}^{\text{BS}}(\mathcal{A}) \doteq \frac{1}{n} \sum_{i=1}^n \frac{1}{|\mathcal{C}^{-i}|} \sum_{b \in \mathcal{C}^{-i}} \ell(y_i, \hat{f}^{*b}(x_i)),$$

where \mathcal{C}^{-i} contains the indices of the bootstraps that do not contain observation (x_i, y_i) .

In order to correct for the optimism of $\hat{\mathcal{R}}^{\text{BS}}(\mathcal{A})$, we combine it with \mathcal{R}^{BS} ,

$$\mathcal{R}^{(0.632)} \doteq 0.368 \hat{\mathcal{R}}^{\text{BS}} + 0.632 \mathcal{R}^{\text{BS}}.$$

Here, a weight of 0.632 is used because that is the probability that a sample (x_i, y_i) appears at least once in a bootstrap sample of size n ,

$$1 - \left(1 - \frac{1}{n}\right)^n \rightarrow 1 - \frac{1}{e} \approx 0.632.$$

This results in a more robust estimator than \mathcal{R}^{BS} .

7.2 Uncertainty in linear models

Suppose we have n observations in our dataset $\{(x_i, y_i)\}_{i=1}^n$. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$ be the design matrix and output vector. We further assume

$$\mathbf{y} \mid \mathbf{X} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}^*, \sigma^2 \mathbf{I}).$$

As we have seen before, the OLSE is computed as follows,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Thus, we have the following distribution over estimators,

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}^*, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}).$$

An unbiased estimator of σ^2 is

$$\hat{\sigma}^2 = \frac{1}{n-d} \sum_{i=1}^n (\hat{\boldsymbol{\beta}}^\top \mathbf{x}_i - y_i)^2.$$

This is the variance computed with the unbiased estimator of $\boldsymbol{\beta}^*$.

Thus, we can approximate a $1 - \alpha$ confidence interval for β_j^* by

$$\hat{\beta}_j \pm z_{\alpha/2} \hat{\epsilon}(\hat{\beta}_j),$$

where $z_{\alpha/2} = \Phi^{-1}(\alpha/2)$, Φ is the CDF of the standard Gaussian, and $\epsilon(\hat{\beta}_j)$ is the j -th diagonal element of $\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})^{-1}$.

7.3 Statistical testing

Assume we have a hypothesis set $\mathcal{H} = \{p(\cdot \mid \boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \Theta\}$ and let $\boldsymbol{\theta}^* \in \Theta$ be the (unknown) true parameter. Furthermore, we are given a null hypothesis $H_0 : \boldsymbol{\theta}^* \in \Theta_0$ and an alternative hypothesis $H_1 : \boldsymbol{\theta}^* \in \Theta_1$ for predefined $\Theta_0, \Theta_1 \subseteq \Theta$. Lastly, we are given n samples $x_1, \dots, x_n \sim p(\cdot \mid \boldsymbol{\theta}^*)$ and a test statistic $t : \mathcal{X}^n \rightarrow \mathbb{R}$. The goal is to find a critical value $c \in \mathbb{R}$ for the test t such that

$$\mathbb{P}(t(X_1, \dots, X_n) \geq c \mid \boldsymbol{\theta})$$

is “low” when $\boldsymbol{\theta} \in \Theta_0$ and “high” when $\boldsymbol{\theta} \in \Theta_1$. So, we accept the null hypothesis when the probability of the test statistic being greater than c is low and we reject it if it is high.

We want to minimize the probability of incorrectly choosing H_1 as this is worse than incorrectly choosing H_0 . We can quantify this notion of risk as

$$\alpha_c \doteq \sup_{\boldsymbol{\theta} \in \Theta_0} \mathbb{P}(t(x_1, \dots, x_n) \geq c \mid \boldsymbol{\theta}).$$

Intuitively, this is the maximum probability of choosing H_1 when H_0 holds. Note that $\alpha_c \rightarrow 0$ as $c \rightarrow \infty$. So, $c^* \rightarrow \infty$ minimizes the risk. But, this comes with the problem that we would never accept H_1 .

The solution to this problem is to forget about choosing the optimal critical value a priori and running an experiment to obtain a realization x_1, \dots, x_n . We then run the test with the realization $t(x_1, \dots, x_n)$ and compute the risk of the least risky critical value that would incorrectly reject H_0 ,

$$p = \inf_{c \in \mathbb{R}} \{ \alpha_c \mid t(x_1, \dots, x_n) \geq c \}.$$

One can show that this is the common p -value,

$$p \doteq \sup_{\theta \in \Theta_0} \mathbb{P}(t(X_1, \dots, X_n) \geq t(x_1, \dots, x_n) \mid \theta).$$

Intuitively, this is the probability that the test is higher than our current observation if we run the experiment again. A very small p -value means that such an extreme observed outcome would be very unlikely under the null hypothesis.

Wald test. The Wald test is an example of a statistical test. Let $\hat{\theta}$ be an estimator of θ with standard deviation $\hat{\sigma}$ that is asymptotically normal. Let the null hypothesis H_0 be $\theta = \theta_0$ and alternative hypothesis H_1 be $\theta \neq \theta_0$, then the Wald test statistic is

$$W = \frac{(\hat{\theta} - \theta_0)^2}{\hat{\sigma}^2}.$$

As can be seen in Figure 7.1, the Wald test is very high when the estimator is confident (low $\hat{\sigma}$) and very different from θ_0 .

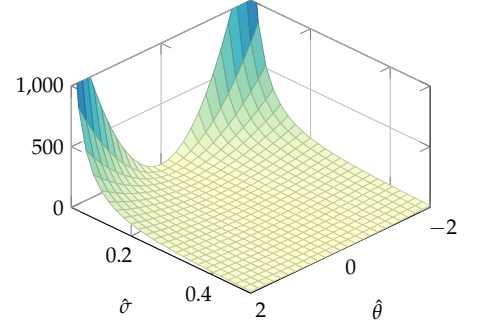


Figure 7.1. Wald test statistic with $\theta_0 = 0$ for different estimators.

7.4 Bayesian neural networks

As we have seen, BLR incorporates Bayesian uncertainty into linear regression. In a similar fashion, BNNs (*Bayesian Neural Networks*) incorporate Bayesian inference into neural networks.

Neural networks. We define a neural network as having L alternating linear layers and element-wise activation functions,

$$\begin{aligned} z_0 &= x \\ z_\ell &= \phi(a_\ell) \\ a_\ell &= W_\ell z_{\ell-1} + b_\ell, \end{aligned}$$

where ϕ is a differentiable activation function and a is its pre-activation. The output of the neural network is then

$$f(x \mid \theta) = z_L, \quad \theta = \{W_\ell, b_\ell\}_{\ell=1}^L.$$

In practice, the weights of the neural network are chosen to minimize the empirical risk,

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \Theta} \hat{\mathcal{R}}(f(\cdot \mid \theta)) \doteq \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i \mid \theta)).$$

In general, the weights are optimized by gradient descent,

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \eta_t \nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}(f(\cdot | \boldsymbol{\theta}_{t-1})).$$

For large datasets, it is too expensive to compute the gradient, so we instead use SGD (*Stochastic Gradient Descent*), which replaces the gradient by an unbiased estimate,

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \eta_t \nabla_{\boldsymbol{\theta}} \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \ell(y_i, f(\mathbf{x}_i | \boldsymbol{\theta}_{t-1})) \right), \quad \mathcal{B} \subseteq [n],$$

where \mathcal{B} contains uniformly sampled sample indices.

Bayesian. The disadvantage of (S)GD is that it yields only a single point estimate of the weights. There is no quantification of uncertainty in this estimate, which leads to overconfidence problems that can result in poor generalization in the presence of domain shifts. We will use a Bayesian approach to alleviate this problem.

We must first define a prior $p(\boldsymbol{\theta})$ over the weights $\boldsymbol{\theta} \in \Theta$,

$$\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

Then, we can define the likelihood function $p(\mathcal{Z} | \boldsymbol{\theta})$ that computes the likelihood of the dataset \mathcal{Z} for a given $\boldsymbol{\theta}$,

$$p(\mathcal{Z} | \boldsymbol{\theta}) = \prod_{(x,y) \in \mathcal{Z}} p(y | \mathbf{x}, \boldsymbol{\theta}),$$

where $p(y | \mathbf{x}, \boldsymbol{\theta})$ is a probability computed by the neural network. Lastly, we use Bayes rule to compute the posterior over weights,

$$p(\boldsymbol{\theta} | \mathcal{Z}) = \frac{p(\boldsymbol{\theta}) p(\mathcal{Z} | \boldsymbol{\theta})}{p(\mathcal{Z})}.$$

The problem is that $p(\mathcal{Z}) = \int p(\boldsymbol{\theta}) p(\mathcal{Z} | \boldsymbol{\theta}) d\boldsymbol{\theta}$ is intractable. The solution is VI (*Variational Inference*)—search a distribution family \mathcal{Q} for the closest distribution to the posterior.⁴ We do this by minimizing the KL divergence between the two,

$$q^* \in \operatorname{argmin}_{q \in \mathcal{Q}} D_{\text{KL}}(q \| p).$$

In our case, we will search the space of isotropic Gaussians. *I.e.*, we search for some mean vector $\boldsymbol{\mu} \in \mathbb{R}^p$ and standard deviation $\sigma > 0$ to form the following distribution,

$$q = \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I}).$$

Thus, we have the following optimization problem,

$$\boldsymbol{\mu}^*, \sigma^* \in \operatorname{argmin}_{\boldsymbol{\mu} \in \mathbb{R}^p, \sigma > 0} D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I}) \| p(\boldsymbol{\theta} | \mathcal{Z})).$$

The Bayesian approach always has three steps: (1) define a prior, (2) define a likelihood, and (3) use Bayes rule to compute the posterior.

⁴ \mathcal{Q} is often called the variational family.

Note that this problem does not have an analytical solution, so we must optimize it with (stochastic) gradient descent. The KL divergence can be rewritten to

$$\begin{aligned} D_{\text{KL}}(q\|p) &\doteq \mathbb{E}_{\theta \sim q} \left[\log \frac{q(\theta)}{p(\theta | \mathcal{Z})} \right] \\ &\propto \mathbb{E}_{\theta \sim q} [\log q(\theta) - \log p(\mathcal{Z} | \theta) - \log p(\theta)]. \end{aligned}$$

Let

$$F(\mu, \sigma, \theta) \doteq \log q(\theta; \mu, \sigma) - \log p(\mathcal{Z} | \theta) - \log p(\theta),$$

then we can derive the gradients to be

$$\begin{aligned} \nabla_{\mu} D_{\text{KL}}(q(\mu, \sigma) \| p(\theta | \mathcal{Z})) &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\nabla_{\theta} F(\mu, \sigma, \theta) + \nabla_{\mu} F(\mu, \sigma, \theta)] \\ \nabla_{\sigma} D_{\text{KL}}(q(\mu, \sigma) \| p(\theta | \mathcal{Z})) &= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [\epsilon^{\top} \nabla_{\theta} F(\mu, \sigma, \theta) + \nabla_{\sigma} F(\mu, \sigma, \theta)]. \end{aligned}$$

We can then approximate the gradients by making a single sample Monte Carlo estimate and apply (S)GD to optimize the weights.

7.5 Information-based transductive learning

We are given a domain \mathcal{X} that contains a safe area $\mathcal{S} \subseteq \mathcal{X}$ and area of interest $\mathcal{A} \subseteq \mathcal{X}$. Furthermore, we have an unknown function $f^* : \mathcal{X} \rightarrow \mathbb{R}$ that we wish to explore within \mathcal{A} . For any point in the safe area $x \in \mathcal{S}$, we can query noisy observations,

$$y_x = f^*(x) + \epsilon_x, \quad \mathbb{E}[\epsilon_x] = 0.$$

However, this is prohibitively expensive in some way, so we want to minimize the amount of queries. So, for every query we make, we want to maximize the amount of information that we receive. Formally, given $n - 1$ previous samples $\mathcal{D}_{n-1} = \{(x_i, y_i)\}_{i=1}^{n-1} \subseteq \mathcal{S} \times \mathbb{R}$, we want to find the next point x_n that will give the most information about f in the area of interest \mathcal{A} .

In ITL (*Information-based Transductive Learning*) [Hübottter et al., 2024], the next point is selected as follows,

$$x_n \in \operatorname{argmax}_{x \in \mathcal{S}} I(f_{\mathcal{A}}; y_x | \mathcal{D}_{n-1}).$$

Intuitively, we are looking for a point $x \in \mathcal{S}$ that maximizes the conditional mutual information between y_x and f restricted to \mathcal{A} . When $f \sim \mathcal{GP}(\mu, k)$ with known mean function μ and kernel k , one can show

$$I(f_{\mathcal{A}}; y_x | \mathcal{D}_{n-1}) = \frac{1}{2} \log \left(\frac{\operatorname{Var}[y_x | \mathcal{D}_{n-1}]}{\operatorname{Var}[y_x | f_{\mathcal{A}}, \mathcal{D}_{n-1}]} \right).$$

We use the reparameterization trick to obtain θ ,

$$\theta = \mu + \sigma \epsilon.$$

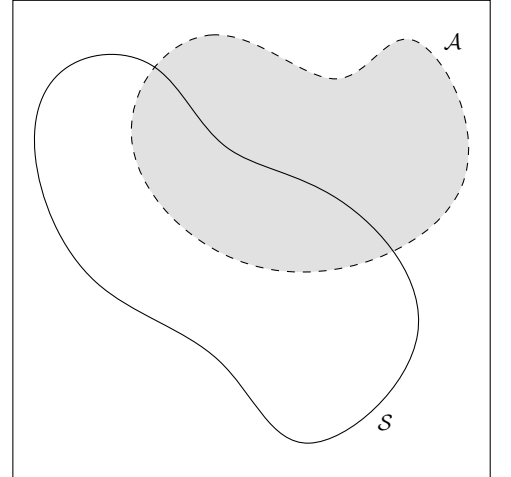


Figure 7.2. Space of active learning.

Proof. Let $X|Y \sim \mathcal{N}(\mu, \sigma^2)$ be conditionally Gaussian, then

$$\begin{aligned} H[X | Y] &= \mathbb{E} \left[-\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(x-\mu)^2}{2\sigma^2} \right) \right) \right] \\ &= \mathbb{E} \left[\log \sqrt{2\pi\sigma^2} - \frac{(x-\mu)^2}{2\sigma^2} \right] \\ &= \frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \mathbb{E}[(x-\mu)^2] \\ &= \frac{1}{2} \log(2\pi) + \frac{1}{2} \log(\sigma^2) - \frac{1}{2} \\ &= \frac{1}{2} (\log(2\pi) + \log \text{Var}[X | Y] - 1). \end{aligned}$$

$$\sigma^2 = \mathbb{E}[(x-\mu)^2].$$

Decomposing mutual information into conditional entropy terms and using that $y_x | f_{\mathcal{A}}, \mathcal{D}_{n-1}$ and $y_x | \mathcal{D}_{n-1}$ are Gaussian yields

$$\begin{aligned} I(f_{\mathcal{A}}, y_x | \mathcal{D}_{n-1}) &= H(y_x | \mathcal{D}_{n-1}) - H(y_x | f_{\mathcal{A}}, \mathcal{D}_{n-1}) \\ &= \frac{1}{2} (\log \text{Var}[y_x | \mathcal{D}_{n-1}] - \log \text{Var}[y_x | f_{\mathcal{A}}, \mathcal{D}_{n-1}]) \\ &= \frac{1}{2} \log \left(\frac{\text{Var}[y_x | \mathcal{D}_{n-1}]}{\text{Var}[y_x | f_{\mathcal{A}}, \mathcal{D}_{n-1}]} \right). \end{aligned}$$

■

Next, we will look at special cases of ITL.

Safe Bayesian optimization. In safe Bayesian optimization, we want to find the maximum of an unknown function f^* , which is done by iteratively choosing points x and observing their realizations $y = f^*(x) + \epsilon_x$. Furthermore, there is an unknown function g^* that defines the safe area,

$$\mathcal{S}^* = \{x \in \mathcal{X} \mid g^*(x) \geq 0\}.$$

We also observe the realizations of this function $z = g^*(x)$. Now, we want to iteratively select points such that we find the maximum of f^* while staying within \mathcal{S}^* .

Assume we are given $n-1$ points again $\{(x_i, y_i, z_i)\}_{i=1}^{n-1}$. We then fit a GP f on $\{(x_i, y_i)\}_{i=1}^{n-1}$, which induces a lower bound function ℓ_n^f and an upper bound function u_n^f , such that $[\ell_n^f(x), u_n^f(x)]$ is the 95% confidence interval of $\mathbb{E}[f(x)]$ for any $x \in \mathcal{X}$. Similarly, we fit a GP g on $\{(x_i, z_i)\}_{i=1}^{n-1}$, which analogously produces two functions ℓ_n^g and u_n^g . These induce the pessimistic and optimistic estimates of the safe set,

$$\mathcal{S}_n \doteq \{x \mid \ell_n^g(x) \geq 0\}, \quad \hat{\mathcal{S}}_n \doteq \{x \mid u_n^g(x) \geq 0\}.$$

Further, we define the area of interest,

$$\mathcal{A}_n \doteq \left\{ x \in \hat{\mathcal{S}}_n \mid u_n^f(x) \geq \max_{x' \in \mathcal{S}_n} \ell_n^f(x') \right\}.$$

Lastly, we can apply ITL using safe set \mathcal{S}_n and area of interest \mathcal{A}_n .

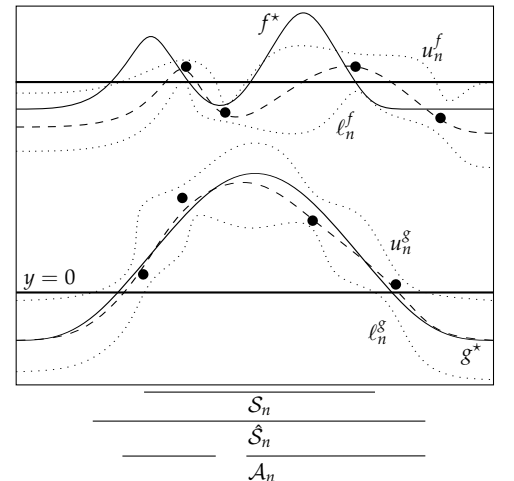


Figure 7.3. Safe Bayesian optimization.

Batch active learning. In batch active learning, we do not select a single point at a time, but we select a batch of points. Formally, we are given an unknown function $f : \mathcal{X} \rightarrow \mathcal{Y}$ that we wish to explore, a population set $X = \{x_1, \dots, x_m\} \subseteq \mathcal{X}$, and a budget $b \leq m$. The goal is to find a subset $L \subseteq X$ such that $|L| = b$ for which we query the oracle $\{f(x) \mid x \in L\}$.

We define $B_\delta(x)$ as the set of the δ -radius ball of points around x ,

$$B_\delta(x) \doteq \{x' \in \mathcal{X} \mid \|x - x'\| \leq \delta\}.$$

We call $B_\delta(x)$ “pure” if f is constant all over $B_\delta(x)$, *i.e.*, every point in a ball around x is classified as the same class. Further, $C(L, \delta)$ is the union of all balls of the chosen subset of points $L \subseteq X$,

$$C(L, \delta) \doteq \bigcup_{x \in L} B_\delta(x).$$

Further, we define

$$\begin{aligned} C_r(L, \delta) &\doteq \{x \in C(L, \delta) \mid \hat{f}(x) = f(x)\} \\ C_w(L, \delta) &\doteq \{x \in C(L, \delta) \mid \hat{f}(x) \neq f(x)\}, \end{aligned}$$

where \hat{f} is a classifier of the space \mathcal{X} . Note $C(L, \delta) = C_r(L, \delta) \cup C_w(L, \delta)$. Lastly, we define $\pi(\delta)$ to be the total probability density of points that have an impure δ -radius ball,

$$\pi(\delta) \doteq \mathbb{P}(\{x \in \mathcal{X} \mid B_\delta(x) \text{ is not pure}\}).$$

Note that π increases with δ . Suppose we know $\mathcal{Z} = \{(x, f(x)) \mid x \in L\}$, then denote the fitted 1-nearest neighbor classifier on this dataset as \hat{f} —the dashed line in Figure 7.4 is its decision boundary.

We have $C_w(L, \delta) \subseteq \{x \in \mathcal{X} \mid B_\delta(x) \text{ is not pure}\}$ (easily verifiable using Figures 7.4 and 7.5). So, we have

$$\mathbb{P}(C_w(L, \delta)) \leq \pi(\delta).$$

Furthermore, looking at Figure 7.4, we can easily see the following

$$\{x \in \mathcal{X} \mid \hat{f}(x) \neq f(x)\} \subseteq C_r^c \cup C_w \subseteq C^c \cup \{x \mid B_\delta(x) \text{ is not pure}\},$$

where A^c denotes the complement of a set A . As a result,

$$\begin{aligned} \mathcal{R}(\hat{f}) &\doteq \mathbb{E}[\mathbb{1}\{\hat{f}(x) \neq f(x)\}] \\ &= \mathbb{P}(\hat{f}(x) \neq f(x)) \\ &\leq 1 - \mathbb{P}(C(L, \delta)) + \pi(\delta). \end{aligned}$$

We want to choose L and δ that minimize $\mathcal{R}(\hat{f})$. This is intractable, so we minimize the upper bound instead. This is done by picking δ first and then choosing L that maximizes $\mathbb{P}(C(L, \delta))$. *I.e.*, we want to solve

$$L^* \in \operatorname{argmax}_{L \subseteq X, |L|=b} \mathbb{P}\left(\bigcup_{x \in L} B_\delta(x)\right).$$

This is the gray area in Figure 7.4.

This is the black area in Figure 7.4.

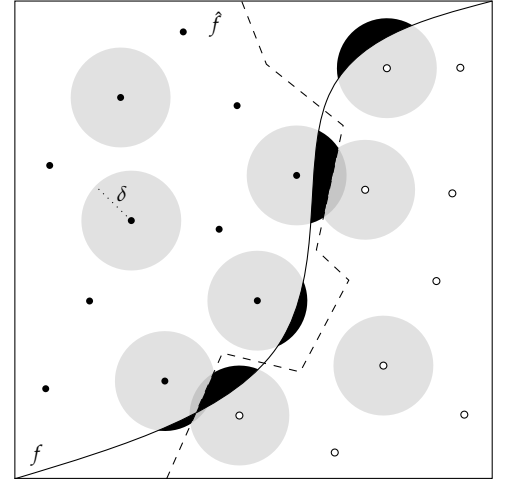


Figure 7.4. Batch active learning. The dashed line is the fitted 1-nearest neighbor classifier on this dataset. Here, $C_w(L, \delta)$ denotes the black area and $C_r(L, \delta)$ denotes the gray area.

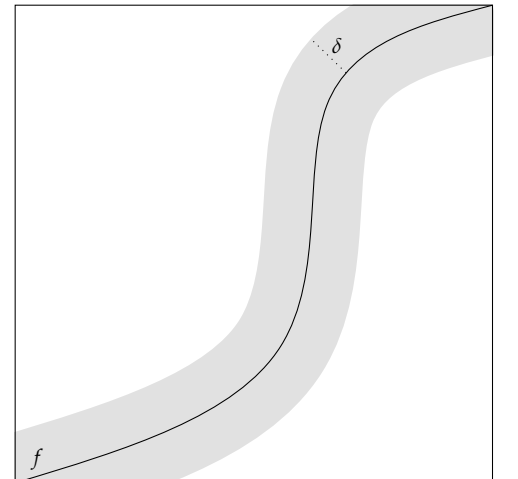


Figure 7.5. $\pi(\delta)$ is the probability density of the marked area, which is the space of impure data points under a given δ .

This has two problems—the distribution of \mathbf{x} is unknown and this problem is NP hard. We address the first problem by approximating the distribution with the empirical distribution induced by X . As a result, we get the following optimization problem,

$$L^* \in \operatorname{argmax}_{L \subseteq X, |L|=b} \frac{1}{|X|} |\{\mathbf{x}' \in X \mid \|\mathbf{x}' - \mathbf{x}\| \leq \delta, \exists \mathbf{x} \in L\}|.$$

We circumvent the second problem by greedily picking points with the most other points in its ball that are not within the ball of previously picked points.

8 Convex optimization

In standard form, an optimization problem looks as follows,

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && g_i(x) = 0, \quad \forall i \in [n] \\ & && h_j(x) \leq 0, \quad \forall j \in [m], \end{aligned}$$

This is called a convex program if f and all h_j are convex and all g_i are affine. We will focus on convex programming. The constraints define a convex set $\mathcal{C} \subseteq \mathcal{X}$, where \mathcal{X} is the intersection of domains of f and all constraints. For our purposes, let $\mathcal{X} = \mathbb{R}^d$. The aim is to compute

$$\operatorname{argmin}_{x \in \mathcal{C}} f(x).$$

Definition 8.1 (Lagrangian dual function). For an optimization problem, the Lagrangian $\mathcal{L} : \mathcal{X} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is defined as

$$\mathcal{L}(x, \lambda, \alpha) \doteq f(x) + \sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \alpha_j h_j(x).$$

The scalars λ_i and α_j are called Lagrange multipliers. The Lagrange dual function $\theta : \mathbb{R}^n \rightarrow \mathbb{R}^m \rightarrow \mathbb{R} \cup \{-\infty\}$ is defined as

$$\theta(\lambda, \alpha) \doteq \inf_{x \in \mathcal{X}} \mathcal{L}(x, \lambda, \alpha).$$

Lemma 8.2 (Weak duality). Let $x \in \mathcal{C}$ (so it satisfies all constraints), $\lambda \in \mathbb{R}^n$, $\alpha \in \mathbb{R}^m$ with $\alpha \geq \mathbf{0}$, and θ be the Lagrange dual function, then

$$\theta(\lambda, \alpha) \leq f(x).$$

The dual Lagrange function lower bounds the primal solution.

Proof. Define x , λ , and α as in Lemma 8.2, then we have

$$\begin{aligned} \theta(\lambda, \alpha) & \doteq \inf_{x \in \mathcal{X}} \mathcal{L}(x, \lambda, \alpha) \\ & \leq \mathcal{L}(x, \lambda, \alpha) \\ & = f(x) + \sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \alpha_j h_j(x) \\ & \leq f(x). \end{aligned}$$

We have $g_i(x) = 0$ and $\alpha_j h_j(x) \leq 0$.

This concludes the proof. ■

As shown in Lemma 8.2, the Lagrange dual function θ lower bounds the primal f for all $x \in \mathcal{C}$, λ , and $\alpha \geq \mathbf{0}$. We are interested in the maximum

lower bound, which results in the dual program,

$$\begin{aligned} & \text{maximize} && \theta(\lambda, \alpha) \\ & \text{subject to} && \alpha_i \geq 0, \quad \forall i \in [n]. \end{aligned}$$

This is a convex program, even if the primal is not convex.

Corollary.

$$\max_{\lambda, \alpha \geq 0} \theta(\lambda, \alpha) \leq \min_{x \in \mathcal{C}} f(x).$$

Lemma 8.3 (Sufficient condition for strong duality). If Slater's condition holds (there exists an $x \in \mathcal{C}$ such that $h_j(x) < 0$ for all $j \in [m]$), then we have strong duality,

$$\max_{\lambda, \alpha \geq 0} \theta(\lambda, \alpha) = \min_{x \in \mathcal{X}} f(x).$$

For most problems, a Slater point exists. So, most of the time, strong duality holds. If all g_i and h_j are differentiable, then the KKT (*Karush-Kuhn-Tucker*) conditions provide necessary conditions for strong duality.

Lemma 8.4 (KKT necessary conditions for strong duality). Let x^* and (λ^*, α^*) be feasible solutions to the primal and dual problems. Further, assume that strong duality holds. Then,

$$\alpha_j^* h_j(x^*) = 0, \quad \forall j \in [m]$$

Complementary slackness.

$$\nabla_x \mathcal{L}(x^*, \lambda^*, \alpha^*) = \nabla f(x^*) + \sum_{i=1}^n \lambda_i^* \nabla g_i(x^*) + \sum_{j=1}^m \alpha_j^* \nabla h_j(x^*) = \mathbf{0}.$$

Vanishing Lagrangian gradient.

The second condition can also be written as the following if \mathcal{L} is finite,

$$x^* \in \operatorname{argmin}_{x \in \mathcal{X}} \mathcal{L}(x, \lambda^*, \alpha^*).$$

Proof. Define x^* , λ^* , and α^* as in Lemma 8.4, then we get the Master equation,

$$\begin{aligned} f(x^*) &= \theta(\lambda^*, \alpha^*) \\ &\doteq \inf_{x \in \mathcal{X}} \mathcal{L}(x, \lambda^*, \alpha^*) \\ &\leq \mathcal{L}(x^*, \lambda^*, \alpha^*) \\ &\doteq f(x^*) + \sum_{i=1}^n \lambda_i^* g_i(x^*) + \sum_{j=1}^m \alpha_j^* h_j(x^*) \\ &\leq f(x^*). \end{aligned}$$

Strong duality.

We have $g_i(x) = 0$ and $\alpha_j h_j(x) \leq 0$.

Thus, all inequalities become equalities. The first condition comes from the second-to-last and last line, which results in $\alpha_j^* h_j(x^*) = 0$ for all

$j \in [m]$. The second condition comes from the second and third line, where there must be vanishing gradient, because \mathbf{x}^* is the minimizer of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*)$. ■

Lemma 8.5 (KKT sufficient conditions for strong duality). Let the primal be a convex program. Further, let \mathbf{x}^* and $(\boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*)$ be feasible solutions to the primal and dual, respectively. Lastly, let the KKT conditions hold,

$$\begin{aligned} \alpha_j^* h_j(\mathbf{x}^*) &= 0, \quad \forall j \in [m] \\ \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*) &= \mathbf{0}. \end{aligned}$$

Then, strong duality holds.

Proof.

$$\begin{aligned} f(\mathbf{x}^*) &= f(\mathbf{x}^*) + \sum_{i=1}^n \lambda_i g_i(\mathbf{x}^*) + \sum_{j=1}^m \alpha_j h_j(\mathbf{x}^*) \\ &= \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*) \\ &= \inf_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*) \\ &\doteq \theta(\boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*). \end{aligned}$$

\mathbf{x}^* is feasible and we have complementary slackness. ■

\mathcal{L} is convex, because convexity is preserved under summation and positive scaling. Also, affine functions remain convex under negative scaling. Furthermore, we have vanishing gradient, so \mathbf{x}^* is a global minimizer of $\mathcal{L}(\cdot, \boldsymbol{\lambda}^*, \boldsymbol{\alpha}^*)$. ■

8.1 Support vector machine

In SVMs (*Support Vector Machines*), we are given a dataset with binary labels and we want to find the linear separator with the maximum margin. Formally, we are given data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \subseteq \mathbb{R}^d \times \{-1, +1\}$ and we model the linear separator as

$$f[\mathbf{w}, b](\mathbf{x}) \doteq \begin{cases} +1 & \mathbf{w}^\top \mathbf{x} + b > 0 \\ -1 & \mathbf{w}^\top \mathbf{x} + b < 0. \end{cases}$$

We want to model this as a convex program, so we need to get rid of the strict inequality and the cases. We are given the dataset, so we can use information of y to reverse the inequality if necessary. Furthermore, we introduce a small scalar $\epsilon > 0$ to get a non-strict inequality,

$$f[\mathbf{w}, b](\mathbf{x}, y) = y(\mathbf{w}^\top \mathbf{x} + b) \geq \epsilon > 0.$$

Now, we have the following optimization problem,

$$\begin{aligned} &\text{maximize} && m(\mathbf{w}, b) \\ &\text{subject to} && y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq \epsilon, \quad \forall i \in [n], \end{aligned}$$

where $m : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$ is the margin function. The margin function is the minimum distance of any point to the linear separator. Let \mathbf{x}^+ and

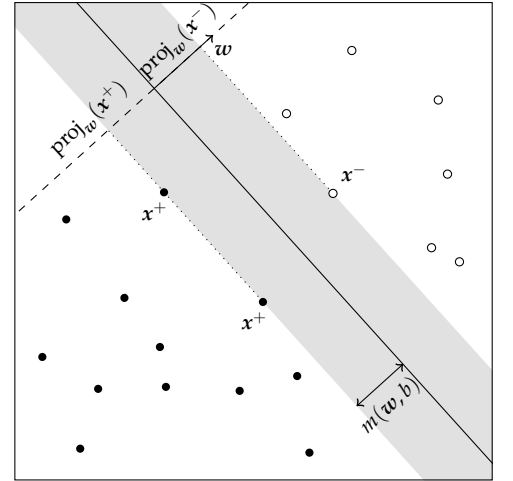


Figure 8.1. Support vector machines.

x^- be the positive and negative support vectors, where a support vector is a point that touches the margin—see Figure 8.1. Then, we have the following definition of the margin,

$$\begin{aligned} 2 \cdot m(w, b) &\doteq \|\text{proj}_w(x^+) - \text{proj}_w(x^-)\| \\ &= \left\| \frac{w^\top x^+}{\|w\|^2} w - \frac{w^\top x^-}{\|w\|^2} w \right\| \\ &= \frac{|w^\top x^+ - w^\top x^-|}{\|w\|^2} \|w\| \\ &= \frac{|w^\top x^+ - w^\top x^-|}{\|w\|} \\ &= \left| \bar{w}^\top (x^+ - x^-) \right|, \end{aligned}$$

See Figure 8.1.

$$\text{proj}_b(a) = \frac{a^\top b}{\|b\|^2} a.$$

where \bar{w} is a normalized w . Note that the margin does not depend on the norm of w or the intercept b —it only depends on the direction of w . As a result, there are an infinite number of solutions, which results in an ill-posed problem. However, only one solution (w, b) will satisfy

$$w^\top x^+ + b = 1, \quad w^\top x^- + b = -1.$$

Under this constraint, the margin can be computed as follows,

$$2 \cdot m(w, b) = \frac{2}{\|w\|},$$

because $|w^\top x^+ - w^\top x^-| = |(1 - b) - (-1 - b)| = 2$. Thus, we get the following convex program,

$$\begin{aligned} &\text{minimize} \quad \frac{1}{2} \|w\|^2 \\ &\text{subject to} \quad 1 - y_i(w^\top x_i + b) \leq 0, \quad \forall i \in [n]. \end{aligned}$$

Minimizing $1/f(x)$ is equivalent to maximizing $f(x)$.

It is easy to see that this convex program has a Slater point. Let (w, b) be a solution to the convex program, so

$$y_i(w^\top x_i + b) \geq 1, \quad \forall i \in [n].$$

Then, $(\gamma w, \gamma b)$ for $\gamma > 1$ is a Slater point, because for all $i \in [n]$, we have

$$y_i(\gamma w^\top x_i + \gamma b) = \gamma y_i(w^\top x_i + b) > y_i(w^\top x_i + b) \geq 1.$$

Hence, strong duality holds and we can use the KKT conditions to find the minimizer x^* .

This convex program has the following Lagrangian,

$$\begin{aligned} \mathcal{L}(w, b, \alpha) &= \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i(w^\top x_i + b)), \quad \alpha_i \geq 0 \\ &= \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i - w^\top \sum_{i=1}^n \alpha_i y_i x_i - b \sum_{i=1}^n \alpha_i y_i. \end{aligned}$$

This function has the following gradients,

$$\begin{aligned}\frac{\partial \mathcal{L}(w, b, \alpha)}{\partial w} &= w - \sum_{i=1}^n \alpha_i y_i x_i \\ \frac{\partial \mathcal{L}(w, b, \alpha)}{\partial b} &= - \sum_{i=1}^n \alpha_i y_i.\end{aligned}$$

From the vanishing Lagrangian gradient KKT condition, we know

$$\frac{\partial \mathcal{L}(w^*, b^*, \alpha)}{\partial w} = 0 \implies w^* = \sum_{i=1}^n \alpha_i y_i x_i.$$

Furthermore, we get the following constraint on α_i ,

$$\frac{\partial \mathcal{L}(w^*, b^*, \alpha)}{\partial b} = 0 \implies \sum_{i=1}^n \alpha_i y_i = 0.$$

Thus, we have the following dual program,

$$\begin{aligned}\text{maximize} \quad & \theta(\alpha) \doteq \mathcal{L}(w^*, b^*, \alpha) \\ \text{subject to} \quad & \alpha_i \geq 0, \quad \forall i \in [n] \\ & \sum_{i=1}^n \alpha_i y_i = 0.\end{aligned}$$

We can simplify the Lagrangian by making use of the found w^* ,

$$\begin{aligned}\theta(\alpha) &\doteq \mathcal{L}(w^*, b^*, \alpha) \\ &= \frac{1}{2} \|w^*\|^2 + \sum_{i=1}^n \alpha_i - \|w^*\|^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \|w^*\|^2 \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^\top x_j.\end{aligned}$$

Thus, the dual program is a quadratic optimization problem on a simplex, which is easy to optimize. Once we have solved the dual program for α^* , we can use this to compute w^* ,

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i.$$

Due to complementary slackness, the following holds,

$$\alpha_i^* \left(1 - y_i (w_*^\top x_i + b^*) \right) = 0.$$

So, $\alpha_i^* = 0$ often, which makes α^* sparse. In fact, $1 - y_i (w_*^\top x_i + b^*) = 0$ only if x_i is a support vector. Thus, the solution w^* is a sparse linear combination of support vectors.

We can compute the intercept by finding support vectors,

$$x^+ \in \operatorname{argmin}_{x_i: y_i=+1} w^\top x_i, \quad x^- \in \operatorname{argmax}_{x_i: y_i=-1} w^\top x_i.$$

Further, remember the constraints $\mathbf{w}_*^\top \mathbf{x}^+ + b = 1$ and $\mathbf{w}_*^\top \mathbf{x}^- + b = -1$, from which we can derive the optimal intercept,

$$b^* = -\frac{1}{2}(\mathbf{w}_*^\top \mathbf{x}^+ + \mathbf{w}_*^\top \mathbf{x}^-).$$

8.2 Soft-margin SVM

The SVM only works if the data is linearly separable, but this does not hold in general. We can work around this by introducing a slackness parameter to each data point, which gives the following convex program,

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i \in [n] \\ & \xi_i \geq 0, \quad i \in [n]. \end{aligned}$$

Here, we introduced slackness parameters that we wish to be as small as possible with some weight C . This allows for incorrect classifications if it is justified. For example, the slack of an outlier might be very high. The hyperparameter C has an important role in balancing the size of the margin with the number of neglected samples. As C increases, the margin will get more narrow, but fewer points will be neglected.

The optimal slackness parameters can be computed as follows,

$$\xi_i^* = \max\{0, 1 - y_i (\mathbf{w}_*^\top \mathbf{x}_i + b^*)\}.$$

And, α_i is further constrained to be less than C .

8.3 Kernelization

The data points might be separated by a non-linear separator. In this case, we can use a different feature space by employing a feature function ϕ . This results in the following solution \mathbf{w}^* ,

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \phi(\mathbf{x}_i).$$

And, we get the following Lagrangian,

$$\begin{aligned} \mathcal{L}(\mathbf{w}^*, \boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j). \end{aligned}$$

So, we do not even need the feature function to compute the Lagrangian—we only need a kernel. Furthermore, we can use the kernel for classification (without loss of generalization, assume $b^* = 0$),

$$\mathbf{w}_*^\top \phi(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}).$$

8.4 Multi-class SVM

We can generalize the notion of a margin to multi-class problems by introducing weights w_z for each class $z \in [k]$. Then, we define the margin as the maximum $m \in \mathbb{R}$ such that

$$m \leq \left(w_{z_i}^\top y_i + b_{z_i} \right) - \max_{z \neq z_i} \left\{ w_z^\top y_i + b_z \right\}, \quad \forall i \in [n].$$

Intuitively, it is the closest distance of any point to a decision boundary. We then define the following optimization problem, which we can solve

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|w\|^2 = \frac{1}{2} \sum_{z=1}^k \|w_z\|^2 \\ & \text{subject to} \quad \left(w_{z_i}^\top y_i + b_{z_i} \right) - \max_{z \neq z_i} \left\{ w_z^\top y_i + b_z \right\} \geq 1, \quad \forall i \in [n]. \end{aligned}$$

For data that is not linearly separable, we can again introduce slack,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad \left(w_{z_i}^\top y_i + b_{z_i} \right) - \max_{z \neq z_i} \left\{ w_z^\top y_i + b_z \right\} \geq 1 - \xi_i, \quad \forall i \in [n] \\ & \quad \xi_i \geq 0, \quad \forall i \in [n]. \end{aligned}$$

8.5 Structural SVM

SVMs can also be used to predict complex objects by making use of its structure. A naive way of achieving this is to give each unique structure its own class and doing multi-class SVM. However, this would result in exponentially (or infinitely) many classes, which is intractable. Thus, the structural SVM needs the following things,

1. Compact representation of the output space;
2. Algorithm that allows for efficient prediction—we cannot enumerate every possible structure;
3. Notion of a prediction error—some structures are more similar than others;
4. Efficient training algorithm that has a runtime complexity that is sub-linear in the number of structures.

The first can be solved by defining a joint feature map $\psi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ that combines properties of inputs and outputs.⁵ Then, we can define the following scoring function,

$$f_w(x, y) = w^\top \psi(x, y).$$

We use this function to perform classification,

$$c(x) = \operatorname{argmax}_{y \in \mathcal{Y}} f_w(x, y).$$

⁵ Notice that the dimensionality of this feature map is independent of the number of structures $|\mathcal{Y}|$.

We would need to construct an algorithm to efficiently compute this argmax. In general, we cannot enumerate over all structures, so we would need to make use of the structure of the data to somehow compute this. For example, the output space might be decomposable into independent parts, such that $\mathcal{Y} = \mathcal{Y}_1 \times \cdots \times \mathcal{Y}_m$.

Moreover, we define the margin as the maximum m that satisfies

$$\mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}_i) - \max_{\mathbf{y} \neq \mathbf{y}_i} \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}) \geq m, \quad \forall i \in [n].$$

Thus, we have the following optimization problem,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to} \quad \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}_i) - \max_{\mathbf{y} \neq \mathbf{y}_i} \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}) \geq 1, \quad \forall i \in [n]. \end{aligned}$$

Here, again the maximum would need a specialized algorithm.

To quantify errors, we introduce a loss function $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, where $\Delta(\mathbf{y}, \mathbf{y}')$ is the loss of predicting \mathbf{y}' when the correct output is \mathbf{y} . Reformulating the optimization problem with slack results in the following,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}_i) - \max_{\mathbf{y} \neq \mathbf{y}_i} \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}) \geq \Delta(\mathbf{y}, \mathbf{y}_i) - \xi_i, \quad \forall i \in [n], \mathbf{y} \neq \mathbf{y}_i. \end{aligned}$$

Or, equivalently with (significantly) less constraints,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{subject to} \quad \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}_i) - \max_{\mathbf{y} \neq \mathbf{y}_i} \left\{ \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}_i) \right\} \geq -\xi_i, \quad \forall i \in [n]. \end{aligned}$$

In conclusion, we need to design the following four elements,

1. Joint feature map $\psi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$;
2. Loss function $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$;
3. Algorithm during training to compute

$$\tilde{\mathbf{y}} \in \operatorname{argmax}_{\mathbf{y} \neq \mathbf{y}_i} \left\{ \mathbf{w}^\top \psi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}_i) \right\};$$

4. Inference algorithm to compute

$$\hat{\mathbf{y}} \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \mathbf{w}^\top \psi(\mathbf{x}, \mathbf{y}).$$

9 Ensembles

The idea of ensembling is simple: create a highly accurate estimator by combining many relatively weak and inaccurate estimators. The advantage over large models is that weak models are easy to train—they only need to be better than random. In general, we have a set of B estimators,

$$\hat{f}_1, \hat{f}_2, \dots, \hat{f}_B.$$

We could take the final estimator to be a convex combination of the estimators,

$$\hat{f}(\mathbf{x}) = \sum_{b=1}^B p_b \hat{f}_b(\mathbf{x}), \quad \sum_{b=1}^B p_b = 1, p_b \geq 0.$$

The bias of this estimator is the same convex combination of the biases of the set of estimators,

$$\begin{aligned} \text{bias}(\hat{f}) &\doteq \mathbb{E}[\hat{f}(X)] - \mathbb{E}[Y | X] \\ &= \sum_{b=1}^B p_b \mathbb{E}[\hat{f}_b(X)] - \mathbb{E}[Y | X] \\ &= \sum_{b=1}^B p_b \left(\mathbb{E}[\hat{f}_b(X)] - \mathbb{E}[Y | X] \right) \\ &= \sum_{b=1}^B p_b \text{bias}(\hat{f}_b). \end{aligned}$$

Thus, unbiased estimators remain unbiased after averaging. We can also look at the variance of this estimator,

$$\begin{aligned} \text{Var}[\hat{f}] &\doteq \mathbb{E} \left[\left(\hat{f}(X) - \mathbb{E}[\hat{f}(X)] \right)^2 \right] \\ &= \mathbb{E} \left[\left(\sum_{b=1}^B p_b \left(\hat{f}_b(X) - \mathbb{E}[\hat{f}_b(X)] \right) \right)^2 \right] \\ &= \frac{1}{B^2} \sum_{b=1}^B \text{Var}[\hat{f}_b] + \frac{1}{B^2} \sum_{b=1}^B \sum_{b' \neq b}^B \text{Cov}[\hat{f}_b, \hat{f}_{b'}]. \end{aligned}$$

Assume $p_b = 1/B$.

Often, the covariances are small ($\text{Cov}[\hat{f}_b, \hat{f}_{b'}] \approx 0$) and the variances similar ($\text{Var}[\hat{f}_b] \approx \sigma^2$), so we gain from averaging because the variance is reduced,

$$\text{Var}[\hat{f}] \approx \frac{\sigma^2}{B}.$$

We will consider techniques for ensembling classifiers, c_1, c_2, \dots, c_B with weights $\alpha_1, \alpha_2, \dots, \alpha_B$. These estimators are composed as follows,

$$\hat{c}(\mathbf{x}) = \text{sgn} \left(\sum_{b=1}^B \alpha_b c_b(\mathbf{x}) \right).$$

It is important to note that this procedure can only succeed if the classifiers are diverse. Otherwise, their covariances are large and we do not gain anything.

Model averaging is very common in machine learning and data science. For example, the Bayesian approach for inference weights parameters according to their posterior,

$$\int p(\boldsymbol{\theta} | \mathbf{X}, \mathbf{y}) p(\mathbf{x}^* | \boldsymbol{\theta}) d\boldsymbol{\theta}.$$

```

1: for  $b \in [B]$  do
2:    $\mathcal{Z}^{*b} = \text{BootstrapSample}(\mathcal{Z})$ 
3:    $\hat{c}_b = \text{TrainClassifier}(\mathcal{Z}^{*b})$ 
4: end for
5: return  $\hat{c}(x) = \text{sgn}\left(\sum_{b=1}^B \hat{c}_b(x)\right)$ 

```

```

1: for  $b \in [B]$  do
2:    $\mathcal{Z}^{*b} = \text{BootstrapSample}(\mathcal{Z})$ 
3:   repeat
4:     Select  $m \approx \sqrt{p}$  features at random
5:     Pick the best feature to split
6:     Split the node into two daughter nodes
7:   until node size  $n_{\min}$  is reached
8: end for
9: return Set of decision trees  $\{\hat{c}_b\}_{b=1}^B$ 

```

Algorithm 2. Bagging with bootstrapped subsets.

Algorithm 3. Random forest algorithm over a dataset \mathcal{Z} , where the feature space is p -dimensional.

9.1 Bagging

Bagging generates diversity in the classifier pool by training on different bootstrapped subsets—see Algorithm 2. Bootstrapped subsets are generated by sampling *i.i.d.* samples from the dataset with replacement.

This works well, because the covariances are small due to using different subsets for training. Furthermore, the variances are similar, because each subsample behaves similarly on average. Finally, the bias is not affected, since the final bias is an average of its estimators' biases.

9.2 Random forests

A random forest is a collection of decision trees, which are used for classification. The idea of the random forest algorithm is to grow a sufficiently deep decision tree to reduce bias, which gives a noisy classifier with high variance—see Algorithm 3. The decision trees are grown by repeatedly splitting nodes based on feature values. An example splitting criterium is splitting to reduce entropy. We can then average to enhance the robustness, since the variance is reduced. Classification is finally done by a majority vote of all decision trees.

9.3 AdaBoost

AdaBoost [Freund and Schapire, 1995] does not use bootstrapped samples, but re-weights the data at each iteration according to the error of previous weak learners. As a result, the weak estimators should be diverse and hence reduce variance. The diversity is due to later classifiers being trained to classify points that previous classifiers found hard with

more weight. The weighting of the data is defined recursively,

$$w_i^{(b+1)} = w_i^{(b)} \exp(\alpha_b \mathbb{1}\{c_b(x_i) \neq y_i\}), \quad w_i^{(1)} = \frac{1}{n}$$

$$\alpha_b \doteq \log\left(\frac{1 - \epsilon_b}{\epsilon_b}\right)$$

$$\epsilon_b \doteq \sum_{i=1}^n \frac{w_i^{(b)}}{\sum_{j=1}^n w_j^{(b)}} \mathbb{1}\{c_b(x_i) \neq y_i\}.$$

We then iteratively fit B classifiers on the full training data using the recursive weights placed on the data points. At inference time, the classifiers are weighted according to their log-odds of error,

$$\hat{c}(x) = \text{sgn}\left(\sum_{b=1}^B \alpha_b c_b(x)\right).$$

Note that AdaBoost is a deterministic algorithm.

Friedman et al. [2000] showed that the minimizer of the exponential loss function⁶ w.r.t. the true distribution is the log-odds of the class probabilities. As such, AdaBoost effectively fits an additive model in base learners, optimizing the exponential function.

Lemma 9.1. The minimizer of $\mathbb{E}[\exp(-yf(x))]$ is

$$f^*(x) = \frac{1}{2} \log \frac{\mathbb{P}(y = 1 | x)}{\mathbb{P}(y = -1 | x)}.$$

Thus, we have the following posteriors,

$$\mathbb{P}(y = 1 | x) = \frac{\exp(f^*(x))}{\exp(-f^*(x)) + \exp(f^*(x))}$$

$$\mathbb{P}(y = -1 | x) = \frac{\exp(-f^*(x))}{\exp(-f^*(x)) + \exp(f^*(x))}$$

Proof. The expected exponential loss function is the following,

$$\mathbb{E}[\exp(-yf(x)) | x] = \mathbb{P}(y = 1 | x) \exp(-f(x)) + \mathbb{P}(y = -1 | x) \exp(f(x))$$

The exponential loss is convex, so the minimizer has vanishing gradient,

$$\frac{\partial \mathbb{E}[\exp(-yf(x)) | x]}{\partial f(x)} = -\mathbb{P}(y = 1 | x) \exp(-f(x)) + \mathbb{P}(y = -1 | x) \exp(f(x))$$

Hence, the minimizer f^* can be derived to be the following (by setting the gradient to zero),

$$\begin{aligned} \mathbb{P}(y = 1 | x) \exp(-f^*(x)) &= \mathbb{P}(y = -1 | x) \exp(f^*(x)) \\ \iff \exp(2 \cdot f^*(x)) &= \frac{\mathbb{P}(y = 1 | x)}{\mathbb{P}(y = -1 | x)} \\ \iff f^*(x) &= \frac{1}{2} \log \frac{\mathbb{P}(y = 1 | x)}{\mathbb{P}(y = -1 | x)}. \end{aligned}$$

This concludes the proof. ■

Log-odds of no error *vs.* error.

Probability of an error, where the normalized weights serve as a prior.

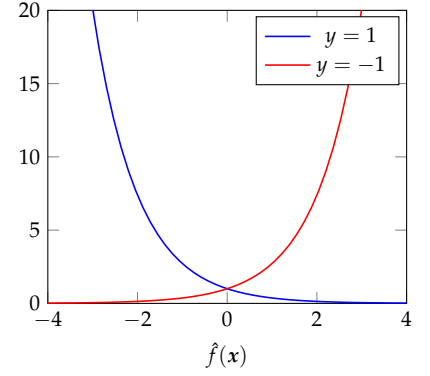


Figure 9.1. Exponential loss function.

⁶ The exponential loss function is $\exp(-yf(x))$.

Theorem 9.2. The discrete AdaBoost algorithm builds an additive logistic regression model via Newton-like updates for minimizing the exponential loss function, $\mathbb{E}[\exp(-yf(\mathbf{x}))]$.

Proof. Let $J(f) = \mathbb{E}[\exp(-yf(\mathbf{x}))]$. Suppose that $f(\mathbf{x})$ is the current estimate and $f(\mathbf{x}) + \alpha c(\mathbf{x})$ is an improvement with $c : \mathcal{X} \rightarrow \{-1, +1\}$ deciding the direction for a small step α . Then, for a fixed α and \mathbf{x} , the second-order Taylor expansion around $c(\mathbf{x}) = 0$ yields

$$\begin{aligned} J(f + \alpha c) &= \mathbb{E}[\exp(-y(f(\mathbf{x}) + \alpha c(\mathbf{x})))] \\ &\approx \mathbb{E}\left[\exp(-yf(\mathbf{x}))\left(1 - y\alpha c(\mathbf{x}) + \frac{1}{2}\alpha^2 y^2 c(\mathbf{x})^2\right)\right] \\ &= \mathbb{E}\left[\exp(-yf(\mathbf{x}))\left(1 - y\alpha c(\mathbf{x}) + \frac{\alpha^2}{2}\right)\right]. \end{aligned}$$

Second-order Taylor approximation around $c(\mathbf{x}) = 0$.

$y^2 = 1$ and $c(\mathbf{x})^2 = 1$.

Minimizing pointwise w.r.t. $c(\mathbf{x}) \in \{-1, 1\}$, we write

$$\begin{aligned} c(\mathbf{x}) &= \underset{c}{\operatorname{argmin}} \mathbb{E}_w\left[1 - y\alpha c(\mathbf{x}) + \frac{\alpha^2}{2} \mid \mathbf{x}\right] \\ &= \underset{c}{\operatorname{argmax}} \mathbb{E}_w[yc(\mathbf{x})]. \end{aligned}$$

where $\mathbb{E}_w[\cdot \mid \mathbf{x}]$ is a weighted conditional expectation, where $w(\mathbf{x}, y) = \exp(-yf(\mathbf{x}))$, and

$$\mathbb{E}_w[g(\mathbf{x}, y) \mid \mathbf{x}] \doteq \frac{\mathbb{E}[w(\mathbf{x}, y)g(\mathbf{x}, y) \mid \mathbf{x}]}{\mathbb{E}[w(\mathbf{x}, y)]}.$$

The solution is

$$c^*(\mathbf{x}) = \begin{cases} 1 & \mathbb{E}_w[y \mid \mathbf{x}] = \mathbb{P}_w(y = 1 \mid \mathbf{x}) - \mathbb{P}_w(y = -1 \mid \mathbf{x}) > 0 \\ -1 & \text{otherwise.} \end{cases}$$

We can directly minimize $J(f + \alpha c)$ without making a Taylor approximation to determine the optimal α^* ,

$$\begin{aligned} \alpha^* &= \underset{\alpha}{\operatorname{argmin}} \mathbb{E}_w[-y\alpha c(\mathbf{x})] \\ &= \frac{1}{2} \log\left(\frac{\mathbb{P}_w(y = 1 \mid \mathbf{x})}{\mathbb{P}_w(y = -1 \mid \mathbf{x})}\right) \\ &= \frac{1}{2} \log\left(\frac{1 - \epsilon_w}{\epsilon_w}\right), \end{aligned}$$

Application of Lemma 9.1.

where $\epsilon_w = \mathbb{E}_w[\mathbb{1}\{y \neq f(\mathbf{x})\}]$.

Combining these steps yields the AdaBoost update,

$$\begin{aligned} f(\mathbf{x}) &\leftarrow f(\mathbf{x}) + \alpha^* c(\mathbf{x}) \\ w(\mathbf{x}, y) &\leftarrow w(\mathbf{x}, y) \exp(-y\alpha^* c(\mathbf{x})) = w(\mathbf{x}, y) \exp(\alpha^* \mathbb{1}\{c(\mathbf{x}) \neq y\}). \end{aligned}$$

This concludes the proof that the AdaBoost algorithm minimizes the exponential loss function by iteratively adding weighted models via second-order optimization. ■

10 Stable diffusion

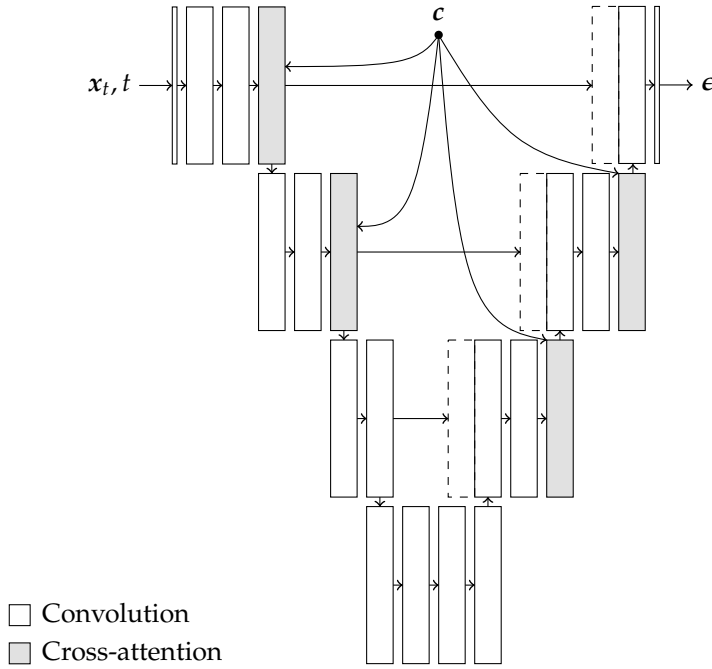


Figure 10.1. Architecture of the stable diffusion denoising model.

10.1 Diffusion models

An SDE (*Stochastic Differential Equation*) is a differential equation in which one or more terms is a stochastic process, resulting in a solution that is also stochastic. Typically, the SDE of a diffusion process is of the following form,

$$dx_t = \mu(x_t, t)dt + \sigma(x_t, t)d\omega_t,$$

where ω_t is a Wiener process (or Brownian motion).⁷ This equation tells us that the change in x_t is driven by a deterministic factor $\mu(x_t, t)$ and a stochastic factor $\sigma(x_t, t)d\omega_t$. Note that $\mu(\cdot, \cdot)$ and $\sigma(\cdot, \cdot)$ induce a probability distribution over time, p_t of x_t .

⁷ Wiener process: $d\omega_t \sim \mathcal{N}(0, dtI)$.

Anderson [1982] showed that the reverse SDE of the diffusion process can be computed as follows,

$$dx_t = \left[\mu(x_t, t) - \sigma^2(x_t, t) \nabla_x \log p_t(x_t) \right] dt + \sigma(x_t, t) d\bar{\omega}_t,$$

where $\bar{\omega}_t$ is a standard Wiener process when time flows backwards from T to 0, and dt is an infinitesimal negative timestep. Thus, in order to do the backward process, we need to learn the score function $\nabla_x \log p_t(x)$ somehow.

Using the DDPM scheduler, diffusion models have the following forward process,

$$x_{t+1} = \sqrt{1 - \beta_t} x_t + \sqrt{\beta_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I).$$

Conditioned on x_t , we can reconstruct x_{t-1} as follows with a predicted noise ϵ_θ ,

$$x_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sqrt{1 - \alpha_t} z, \quad z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{\tau=1}^t \alpha_\tau$. Thus, we need to learn this function.

It can be shown that a diffusion model with the DDPM scheduler is an approximation of a discretization of the following SDE,

$$dx_t = -\frac{1}{2}\beta_t x_t dt + \sqrt{\beta_t} d\omega_t.$$

The reverse process is thus given by the following reverse SDE,

$$dx_t = \left[-\frac{1}{2}\beta_t x_t - \beta_t \nabla_x \log p(x_t) \right] dt + \sqrt{\beta_t} d\bar{\omega}_t.$$

In practice, we train a diffusion model by randomly sampling $x_0 \sim p_0, t \sim \text{Unif}([T]), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and performing a gradient step on the following loss function,

$$\ell = \|\epsilon - \epsilon_\theta(x_t)\|^2.$$

We can sample by iteratively denoising using Equation (1), starting from $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. It can be shown that $\epsilon_t \approx \nabla_x \log p_t(x_t)$.

10.2 U-net

U-nets [Ronneberger et al., 2015] are models used for image-to-image translation tasks. In the case of diffusion models, we have such a task, where we get x_t as input and want to predict ϵ . This framework is used to model ϵ_θ . U-nets work by downsampling the input in stages and then upsampling back to the original space in the same stages. At every step of upsampling, the output of the corresponding downsampling step is concatenated to its input. As such, each decoding step receives low-level and high-level information as input.

10.3 Latent diffusion models

Stable diffusion [Rombach et al., 2022] performs diffusion modeling in the latent space of a pretrained VAE [Kingma, 2013]. During training, we thus first map the input image $x_0 \in \mathbb{R}^d$ to its latent encoding,

$$z_0 = \mathcal{E}(x_0), \quad z_0 \in \mathbb{R}^{d'}, \quad d' \ll d.$$

Then, we use the same loss function as above, where we sample a random timestep and noise,

$$\ell = \|\epsilon - \epsilon_\theta(z_t)\|^2, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d'}).$$

Drift term $\mu(x_t, t) \doteq -\frac{1}{2}\beta_t x_t$; diffusion term $\sigma(t) \doteq \sqrt{\beta_t}$. We can show that this approximates the diffusion model by discretizing,

$$\begin{aligned} x_{t+1} - x_t &= -\frac{1}{2}\beta_t x_t + \sqrt{\beta_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ x_{t+1} &= \left(1 - \frac{1}{2}\beta_t\right) x_t + \sqrt{\beta_t} \epsilon \\ &\approx \sqrt{1 - \beta_t} x_t + \sqrt{\beta_t} \epsilon. \end{aligned}$$

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \text{ where } \bar{\alpha}_t = \prod_{\tau=1}^t (1 - \beta_\tau).$$

Then, during inference, we sample a noise vector in the latent space $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d'})$ and denoise using Equation (1) to get $\tilde{\mathbf{z}}_0$. Lastly, we decode the latent vector back into pixel space,

$$\tilde{\mathbf{x}}_0 = \mathcal{D}(\tilde{\mathbf{z}}_0).$$

The reconstruction term of this VAE is most important, because the model will be searching through the latent space anyway. Stable diffusion uses a latent space that is downsampled 8 times and has 4 channels. This results in significant savings on computation time.

10.4 Text embeddings

In order to perform text-to-image generation, we will need a continuous high-dimensional representation of the input text. For this, we use CLIP [Radford et al., 2021]. CLIP trains image and text transformer models to align text-image pairs. Because of this, they contain more semantic embeddings than other methods of training models to obtain text embeddings. Given a text input sequence of size T_c , the CLIP text model returns a sequence of embeddings,

$$\mathbf{C} \in \mathbb{R}^{T_c \times d_c}.$$

We denote this matrix by \mathbf{C} , because we use it for conditioning.

This sequence is contextualized, because CLIP makes use of self-attention.

10.5 Cross-attention

Now the question becomes how to condition a U-net on the input text sequence $\mathbf{C} \in \mathbb{R}^{T_c \times d_c}$. Stable diffusion [Rombach et al., 2022] does this by making use of cross-attention blocks, which it places at the end of downsampling and upsampling stages of the U-net. It first rearranges the output of the downsampling block into timesteps,

$$\mathbf{X} \in \mathbb{R}^{T \times d}.$$

Then, it computes queries from \mathbf{X} , and keys and values from \mathbf{C} ,

$$\begin{aligned} \mathbf{Q} &= \mathbf{X}\mathbf{W}_Q, & \mathbf{W}_Q &\in \mathbb{R}^{d_k \times d} \\ \mathbf{K} &= \mathbf{C}\mathbf{W}_K, & \mathbf{W}_K &\in \mathbb{R}^{d_k \times d_c} \\ \mathbf{V} &= \mathbf{C}\mathbf{W}_V, & \mathbf{W}_V &\in \mathbb{R}^{d_v \times d_c}. \end{aligned}$$

It uses these to perform the attention mechanism with a residual connection,

$$\mathbf{\Xi} = \mathbf{X} + \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V},$$

where $\mathbf{\Xi}$ denotes the conditioned representation of \mathbf{X} . Note that this architecture is agnostic to the type of the condition. As long as we can embed the conditioning variable, we can condition on it in this way—e.g., we can additionally condition on images [Ye et al., 2023].

11 Non-parametric Bayesian methods

The common beta distribution is defined as follows,

$$\text{Beta}(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad x \in [0, 1], \alpha, \beta > 0.$$

$$B(\alpha, \beta) \doteq \Gamma(\alpha)\Gamma(\beta)/\Gamma(\alpha+\beta), \text{ where } \Gamma(\alpha) = \int_0^\infty \exp(-x)x^{\alpha-1}dx.$$

Intuitively, the Beta distribution is the probability of a Bernoulli process with $p = x$ after observing $\alpha - 1$ successes and $\beta - 1$ failures. The Dirichlet distribution is the multivariate generalization of the beta distribution,

$$\text{Dir}(\mathbf{x}; \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^n x_k^{\alpha_k-1}, \quad \mathbf{x} \in [0, 1]^n, \boldsymbol{\alpha} \in \mathbb{R}_{>0}^n.$$

$$B(\boldsymbol{\alpha}) \doteq \prod_{k=1}^n \Gamma(\alpha_k) / \Gamma(\sum_{k=1}^n \alpha_k).$$

We can model the model generative process with a Dirichlet distribution by sampling the cluster means from a fixed distribution,

$$\mu_1, \dots, \mu_K \stackrel{i.i.d.}{\sim} \mathcal{N}(\mu_0, \sigma_0).$$

Then, we sample the cluster probabilities from the Dirichlet distribution,

$$(\rho_1, \dots, \rho_K) \sim \text{Dir}(\alpha_1, \dots, \alpha_K).$$

We assign the points to their clusters by

$$z_i \sim \text{Categorical}(\rho_1, \dots, \rho_K).$$

Lastly, the coordinates of the data points are sampled from their cluster's distribution,

$$x_i \sim \mathcal{N}(\mu_{z_i}, \sigma).$$

The problem with models like the GMM is that we need to specify the number of classes K beforehand. However, generally, there is no way of knowing the number of classes. Non-parametric Bayesian methods solve this by setting $K = \infty$, meaning that we have infinitely many parameters and we can keep drawing new parameters. However, we cannot draw infinite points from a Dirichlet distribution.

11.1 Dirichlet process

A Dirichlet process $DP(\alpha, H)$ is a distribution over probability distributions on a space Θ , where $\alpha > 0$ is the concentration parameter and H is the base measure on Θ . Specifically, a sample $G \sim DP(\alpha, H)$ is a function $G : \Theta \rightarrow \mathbb{R}_{\geq 0}$ such that $\int_{\Theta} G(\theta) d\theta = 1$. $DP(\alpha, H)$ is characterized by the following property: for every partition (T_1, \dots, T_k) of Θ and $G \sim DP(\alpha, H)$, we have

$$(G(T_1), \dots, G(T_k)) \sim \text{Dir}(\alpha H(T_1), \dots, \alpha H(T_k)).$$

One important observation of the Dirichlet distribution $\text{Dir}(\alpha_1, \dots, \alpha_K)$ is that we can equivalently define its sampling recursively,

$$\begin{aligned}\beta_i &\sim \text{Beta}\left(\alpha_i, \prod_{k=i+1}^K \alpha_k\right) \\ \rho_i &= \beta_i \prod_{j=1}^{i-1} (1 - \beta_j) \\ (\rho_{i+1}, \dots, \rho_K) &\sim \text{Dir}(\alpha_{i+1}, \dots, \alpha_K),\end{aligned}$$

where we have the following base case,

$$\rho_K = 1 - \sum_{k=1}^K \rho_k.$$

This is called stick-breaking. However, we are still limited to finite K . The GEM distribution fixed this by fixing α such that $\beta_i \sim \text{Beta}(1, \alpha)$ for all i . So, we get the following recursion,

$$\begin{aligned}\beta_i &\sim \text{Beta}(1, \alpha) \\ \rho_i &= \beta_i \prod_{j=1}^{i-1} (1 - \beta_j),\end{aligned}$$

where the base case is

$$\rho_K = \beta_K \left(1 - \sum_{i=1}^{K-1} \rho_i\right).$$

Here, we can keep sampling until we are satisfied with K .

If $(\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha)$ and $\theta_k \sim H$ for $k = 1, 2, \dots$, then the following is a sample from $DP(\alpha, H)$,

$$G(\theta) = \sum_{k=1}^{\infty} \rho_k \delta_{\theta_k}(\theta).$$

In particular, samples from a DP are almost surely discrete measures on Θ , supported on a countable set $\{\theta_k\}_{k=1}^{\infty} \subset \Theta$.

11.2 Chinese restaurant process

In order to present a technique for drawing samples from a Dirichlet process, we use a Chinese restaurant metaphor. Here, customers are observations and tables are clusters. When a new customer arrives, he can either join an existing table with a probability proportional to the number of people sitting there or start a new table with probability proportional to $\alpha > 0$. For example, we might have the following partition,

$$\mathcal{P} = \{\{1, 2, 4, 7\}, \{3, 6, 8\}, \{5, 9\}, \{10\}\}.$$

Here, a table is represented by a set of its customers. Using this notation, we can define the probability distribution of a customer joining a table

given the current partition as

$$P(\text{customer } n+1 \text{ joins table } \theta \mid \mathcal{P}) = \begin{cases} \frac{|\theta|}{\alpha+n} & \theta \in \mathcal{P} \\ \frac{\alpha}{\alpha+n} & \text{otherwise.} \end{cases}$$

The above example partition has the following probability, where the number of the customers indicate the order in which they entered the restaurant,

$$P(\mathcal{P}) = \left(\frac{\alpha}{\alpha}\right) \left(\frac{1}{\alpha+1}\right) \left(\frac{\alpha}{\alpha+2}\right) \left(\frac{2}{\alpha+3}\right) \left(\frac{\alpha}{\alpha+4}\right) \cdots$$

An α in the numerator means that a new table was created.

This can be written in a closed form,

$$P(\mathcal{P}) = \alpha^{|\mathcal{P}|} \frac{\alpha!}{(N+\alpha)!} \prod_{\tau \in \mathcal{P}} (|\tau| - 1)!$$

This means that this process is exchangeable, meaning that it is order and labeling independent. The expected number of tables is on the order of $\mathcal{O}(\alpha \log N)$, which means that the expected cluster count is a prior property of a Dirichlet process.

11.3 Dirichlet process mixture model

Let Θ be a set that parametrizes a set of probability distributions and fix a base measure H on Θ . For example, we might have $\Theta = \mathbb{R}$ with $\mu \in \Theta$ corresponding to $\mathcal{N}(\mu, \sigma)$ for some fixed $\sigma > 0$ and $H = \mathcal{N}(\mu_0, \sigma_0)$ for some fixed $\mu_0 \in \mathbb{R}, \sigma_0 > 0$. Using that, we can define the generative process of a DP mixture model, where the probabilities of the clusters are sampled from a GEM distribution,

$$(\rho_1, \rho_2, \dots) \sim \text{GEM}(\alpha).$$

The centers of the clusters are sampled from the base measure,

$$\mu_1, \mu_2, \dots \stackrel{i.i.d.}{\sim} H = \mathcal{N}(\mu_0, \sigma_0).$$

We assign data points to the cluster based on the probabilities,

$$z_i \sim \text{Categorical}(\rho_1, \rho_2, \dots), \quad i \in [n].$$

And, we sample the coordinates of the data points from the cluster distributions,

$$x_i \sim \mathcal{N}(\mu_{z_i}, \sigma), \quad i \in [n].$$

By exchangeability, we can consider any point to be the “last arrived”.

To fit a DPMM, we use Gibbs sampling, where the idea is that we sample each variable in turn, conditioned on the values of all the other

variables in the distribution. Collapsed Gibbs sampling formulation:

$$\begin{aligned}
 p(z_i = k \mid \mathbf{z}^{-i}, \mathbf{x}, \alpha, \boldsymbol{\mu}) & \\
 &\propto p(z_i = k \mid \mathbf{z}^{-i}, \alpha, \boldsymbol{\mu}) p(\mathbf{x} \mid z_i = k, \mathbf{z}^{-i}, \alpha, \boldsymbol{\mu}) \\
 &= p(z_i = k \mid \mathbf{z}^{-i}, \alpha) p(\mathbf{x} \mid z_i = k, \mathbf{z}^{-i}, \boldsymbol{\mu}) \\
 &= p(z_i = k \mid \mathbf{z}^{-i}, \alpha) p(x_i \mid \mathbf{x}^{-i}, z_i = k, \mathbf{z}^{-i}, \boldsymbol{\mu}) p(\mathbf{x}^{-i} \mid \mathbf{z}^{-i}, \boldsymbol{\mu}) \\
 &\propto \underbrace{p(z_i = k \mid \mathbf{z}^{-i}, \alpha)}_{\text{Prior}} \underbrace{p(x_i \mid \mathbf{x}^{-i}, z_i = k, \mathbf{z}^{-i}, \boldsymbol{\mu})}_{\text{Likelihood}}.
 \end{aligned}$$

Bayes rule.

z_i is independent of $\boldsymbol{\mu}$ and \mathbf{x} is independent of α .

Product rule.

$p(\mathbf{x}^{-i} \mid \mathbf{z}^{-i}, \boldsymbol{\mu})$ is a constant factor.

The prior is simple to compute, because we can view i as the last point to enter the dataset—akin to how we can view any customer as the last to enter the Chinese restaurant due to exchangeability. So, we get

$$p(z_i = k \mid \mathbf{z}^{-i}, \alpha) = \begin{cases} \frac{N_k^{-i}}{\alpha + N - 1} & \text{for existing } k \\ \frac{\alpha}{\alpha + N - 1} & \text{otherwise,} \end{cases}$$

where N_k^{-i} is the number of elements in cluster k , excluding i .

The likelihood is conditioned on the cluster k , so we only have to consider points within that cluster. Let \mathbf{x}_k^{-i} be the set of points assigned to cluster k , excluding i . Then,

$$p(x_i \mid \mathbf{x}^{-i}, z_i = k, \mathbf{z}^{-i}, \boldsymbol{\mu}) = \begin{cases} p(x_i \mid \mathbf{x}_k^{-i}, \boldsymbol{\mu}) = \frac{p(x_i, \mathbf{x}_k^{-i} \mid \boldsymbol{\mu})}{p(\mathbf{x}_k^{-i} \mid \boldsymbol{\mu})} & \text{for existing } k \\ p(x_i \mid \boldsymbol{\mu}) & \text{otherwise.} \end{cases}$$

So, the final Gibbs sampler iterates over the points $i \in [N]$ in a random order and samples a new cluster allocation from the following distribution,

$$p(z_i = k \mid \mathbf{z}^{-i}, \mathbf{x}, \alpha, \boldsymbol{\mu}) = \begin{cases} \frac{N_k^{-i}}{\alpha + N - 1} p(x_i \mid \mathbf{x}_k^{-i}, \boldsymbol{\mu}) & \text{for existing } k \\ \frac{\alpha}{\alpha + N - 1} p(x_i \mid \boldsymbol{\mu}) & \text{otherwise.} \end{cases}$$

As a result, we add new clusters with a small probability at each iteration, dependent on α . Furthermore, points are more likely to be added to clusters with many points. Also, the probability of joining a cluster is dependent on the cluster means.

12 Statistical learning theory

12.1 PAC learning

PAC (*Probably Approximately Correct*) learning is a subfield of machine learning that concerns itself with the definition of “learnable” and if something is “learnable”, whether it can be learned by empirically minimizing a loss function. We will be using the following fundamental definitions and notations,

- \mathcal{X} : Data space over which there is some distribution p ;
- \mathcal{C} : Concept class, whose elements decide the labels of \mathcal{X} . We can think of it as a set of functions $c : \mathcal{X} \rightarrow \mathcal{Y}$;
- \mathcal{H} : Hypothesis class, which is a space of solutions that the model can consider, *e.g.*, the weight space of a linear model. Generally, $|\mathcal{H}| \ll |\mathcal{C}|$.

No additional prior knowledge on p is available, which makes this setting differ from Bayesian approaches, which require a prior on \mathcal{X} .

The setting is that nature decides the distribution over \mathcal{X} and its concept class \mathcal{C} . We draw samples from \mathcal{X} according to p and observe its labels from \mathcal{C} —this forms the dataset $\mathcal{Z} \subseteq \mathcal{X} \times \mathcal{Y}$. The learning algorithm receives \mathcal{Z} and \mathcal{H} as input and selects a hypothesis $\hat{c} \in \mathcal{H}$. Lastly, the hypothesis \hat{c} is evaluated on its generalization error,

$$\mathcal{R}(\hat{c}) \doteq \mathbb{P}(\hat{c}(X) \neq c(X)) = \mathbb{E}[\mathbb{1}\{\hat{c}(X) \neq c(X)\}].$$

In practice, the expected risk is approximated using held-out validation data, since we do not have access to the distribution p over X .

Definition 12.1 (Learning algorithm). Let \mathcal{H} and $c \in \mathcal{C}$ be a hypothesis class and a concept. A learning algorithm \mathcal{A} is an algorithm that receives as input a labeled dataset $\mathcal{Z} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ with $y_i = c(x_i)$ for all $i \in [n]$ and outputs a hypothesis $\hat{c} \in \mathcal{H}$.

The objective of the learning algorithm is to minimize the generalization error. However, this is intractable, because we do not have access to p , only its samples. Hence, we approximate it by the empirical error,

$$\hat{c} \in \operatorname{argmin}_{c \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\hat{c}(x_i) \neq c(x_i)\}.$$

Definition 12.2 (Learnable). A learning algorithm \mathcal{A} can learn a concept $c \in \mathcal{C}$ if there exists a polynomial function $\text{poly}(\cdot, \cdot, \cdot)$ such that for any distribution p on \mathcal{X} and $\epsilon, \delta \in (0, 1/2)$, if \mathcal{A} receives as input a sample \mathcal{Z} of size $n \geq \text{poly}(1/\epsilon, 1/\delta, \text{size}(c))$, then \mathcal{A} outputs \hat{c} such that

$$\mathbb{P}(\mathcal{R}(\hat{c}) \leq \epsilon) \geq 1 - \delta.$$

Note that this probability is taken over the randomness of \mathcal{Z} and the potential randomness of \mathcal{A} . The value $\text{size}(c)$ indicates the size of the representation of concept c .

Here, ϵ is an error parameter and δ denotes a confidence value.

Definition 12.3 (PAC learnable). A concept class \mathcal{C} is PAC learnable from a hypothesis class \mathcal{H} if there is an algorithm that can learn any concept in \mathcal{C} .

Definition 12.4 (Efficient PAC learning). If \mathcal{A} runs in time polynomial in only $1/\epsilon$ and $1/\delta$, \mathcal{C} is efficiently PAC learnable.

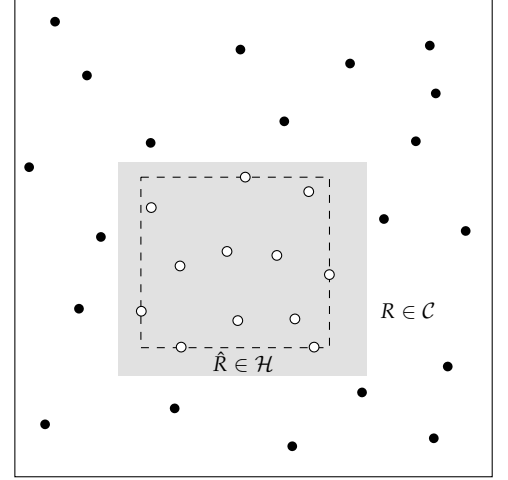


Figure 12.1. The dots, along with their color, represent the dataset $\mathcal{Z} \sim p$. The concept class \mathcal{C} and hypothesis class \mathcal{H} is the set of all rectangles in \mathbb{R}^2 .

12.2 Axis-aligned rectangles

We will now consider the problem of learning an axis-aligned rectangle—see Figure 12.1—and we will prove that this problem is PAC learnable. For this, we need to show that there exists a $\text{poly}(\cdot, \cdot, \cdot)$ such that for any $R \in \mathcal{C}$, p on \mathbb{R}^2 , and $\epsilon, \delta > 0$, the condition in Definition 12.2 holds. In this problem, $\text{size}(R) = 4$. The algorithm we use is taking the tightest rectangle around the points.

Let $\epsilon, \delta > 0$ and $R \in \mathcal{C}$. We now define an event $\hat{R}IG$ (\hat{R} Is Good Enough). Let T_{\uparrow}^{ϵ} be the upper strip of R such that $\mathbb{P}(T_{\uparrow}^{\epsilon}) = \epsilon/4$ —see Figure 12.2.⁸ Analogously define $T_{\downarrow}^{\epsilon}, T_{\leftarrow}^{\epsilon}, T_{\rightarrow}^{\epsilon}$.

The $\hat{R}IG$ event occurs when a hypothesis \hat{R} intersects all four T_{*}^{ϵ} . This only occurs if there is a sample from p within every T_{*}^{ϵ} . We have constructed these spaces, such that the probability of a single sample being within one is $\epsilon/4$. Hence, the probability that a single sample from p misses is $1 - \epsilon/4$, so the probability that n samples miss is $(1 - \epsilon/4)^n$. The probability of the $\hat{R}IG$ event is thus

$$\mathbb{P}(\hat{R}IG) \geq 1 - 4 \cdot \left(1 - \frac{\epsilon}{4}\right)^n \geq 1 - 4 \exp\left(-\frac{n\epsilon}{4}\right).$$

The risk is the density of the difference between the sets R and \hat{R} ,

$$\mathcal{R}(\hat{R}) = \mathbb{P}(\hat{R}(X) \neq R(X)) = \mathbb{P}((R \setminus \hat{R}) \cup (\hat{R} \setminus R)) = \mathbb{P}(R \setminus \hat{R}).$$

Since $\mathbb{P}(R \setminus \hat{R}) \leq \epsilon$ is guaranteed in—but not limited to—the event of

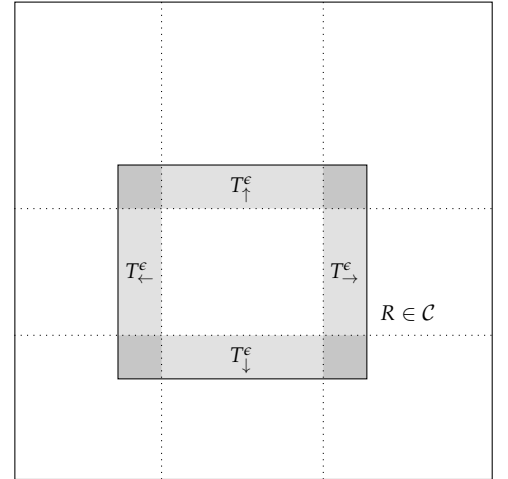


Figure 12.2. The event $\hat{R}IG$ is the event that \hat{R} intersects all four strips.

⁸ $\mathbb{P}(A)$, where A is a set, is the probability of a data point $x \sim p$ being in the set $x \in A$.

$$1 - x \leq \exp(-x).$$

$\hat{R}IG$, we have the following bound on the generalization error,

$$\mathbb{P}(\mathcal{R}(\hat{R}) \leq \epsilon) = \mathbb{P}(\mathbb{P}(R \setminus \hat{R}) \leq \epsilon) \geq \mathbb{P}(\hat{R}IG).$$

As a result,

$$\mathbb{P}(\mathcal{R}(\hat{R}) \leq \epsilon) \geq 1 - 4 \exp\left(-\frac{n\epsilon}{4}\right).$$

For $n \geq \frac{4}{\epsilon} \log\left(\frac{4}{\delta}\right)$, we then have

$$\mathbb{P}(\mathcal{R}(\hat{R}) \leq \epsilon) \geq 1 - \delta.$$

This concludes the proof.

12.3 Stochastic setting

In general, the label of an instance $x \in \mathcal{X}$ is not deterministically determined by the underlying concept—it is modeled by a distribution on $\mathcal{X} \times \mathcal{Y}$. This reflects that two instances with identical features may have different labels. Thus, the definition of a dataset is altered, where it is now $\mathcal{Z} = \{(x_1, y_1), \dots, (x_n, y_n)\}$, where $(x_i, y_i) \sim p$ for all $i \in [n]$. In this setting, the goal is to find a hypothesis $\hat{c} \in \mathcal{H}$ with small generalization error,

$$\mathcal{R}(\hat{c}) = \mathbb{P}_{x,y \sim p}(\hat{c}(x) \neq y) = \mathbb{E}_{x,y \sim p}[\mathbb{1}\{\hat{c}(x) \neq y\}].$$

If the Bayes-optimal classifier is not an element of the concept class \mathcal{C} , then it is impossible to construct an algorithm that finds a $\hat{c} \in \mathcal{H}$ such that $\mathcal{R}(\hat{c}) \leq \epsilon$ for all $\epsilon \in (0, 1/2)$ for a polynomial number of data points. Thus, instead we aim to minimize the generalization error to the best concept,

$$\mathcal{R}(\hat{c}) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) \leq \epsilon.$$

Definition 12.5 (Learnable in the stochastic setting). A learning algorithm \mathcal{A} can learn a concept $c \in \mathcal{C}$ from \mathcal{H} if there is a $\text{poly}(\cdot, \cdot, \cdot)$ such that for any distribution p on $\mathcal{X} \times \mathcal{Y}$ and for any $\epsilon, \delta \in (0, 1/2)$, if \mathcal{A} receives as input a sample \mathcal{Z} of size $n \geq \text{poly}(1/\epsilon, 1/\delta, \text{size}(c))$, then \mathcal{A} outputs $\hat{c} \in \mathcal{H}$ such that

$$\mathbb{P}_{\mathcal{Z} \sim p} \left(\mathcal{R}(\hat{c}) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) \leq \epsilon \right) \geq 1 - \delta.$$

12.4 Vapnik-Chervonenkis

The VC (*Vapnik Chervonenkis*) dimension is defined as the cardinality of the largest set of points that a concept class \mathcal{C} can shatter, which means that there exists a concept in \mathcal{C} for every labeling of the points. Intuitively, this is a measure of the complexity of the concept class.

Theorem 12.6 ([Vapnik and Chervonenkis, 1974]). Assume a finite concept class such that $|\mathcal{C}| < \infty$ and $\mathcal{R}(c^*) = 0$ for the optimal $c^* \in \mathcal{C}$. Further, define the empirical risk minimizer c_n^* as

$$c_n^* \in \{c \in \mathcal{C} \mid \hat{\mathcal{R}}_n(c) = 0\}.$$

Then, for every $n \in \mathbb{N}$ and $\epsilon > 0$,

$$\mathbb{P}(\mathcal{R}(c_n^*) > \epsilon) \leq |\mathcal{C}| \exp(-n\epsilon).$$

And,

$$\mathbb{E}[\mathcal{R}(c_n^*)] \leq \frac{1 + \log |\mathcal{C}|}{n}.$$

Proof. We will only show the first statement,

$$\begin{aligned} \mathbb{P}(\mathcal{R}(c_n^*) > \epsilon) &\leq \mathbb{P}\left(\max_{c \in \mathcal{C}: \hat{\mathcal{R}}_n(c)=0} \mathcal{R}(c) > \epsilon\right) \\ &= \mathbb{P}\left(\left(\max_{c \in \mathcal{C}} \mathbb{1}\{\hat{\mathcal{R}}_n(c) = 0\} \cdot \mathcal{R}(c)\right) > \epsilon\right) \\ &= \mathbb{E}\left[\max_{c \in \mathcal{C}} \mathbb{1}\{\hat{\mathcal{R}}_n(c) = 0\} \cdot \mathbb{1}\{\mathcal{R}(c) > \epsilon\}\right] \\ &\leq \sum_{c \in \mathcal{C}: \mathcal{R}(c) > \epsilon} \mathbb{E}[\mathbb{1}\{\hat{\mathcal{R}}_n(c) = 0\}] \\ &= \sum_{c \in \mathcal{C}: \mathcal{R}(c) > \epsilon} \mathbb{P}(\hat{\mathcal{R}}_n(c) = 0) \\ &\leq |\mathcal{C}|(1 - \epsilon)^n \\ &\leq |\mathcal{C}| \exp(-n\epsilon). \end{aligned}$$

This concludes the proof. ■

Theorem 12.7 (VC inequality).

$$\mathbb{P}\left(\mathcal{R}(c_n^*) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) > \epsilon\right) \leq \mathbb{P}\left(\sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \frac{\epsilon}{2}\right).$$

Proof.

$$\begin{aligned} \mathcal{R}(c_n^*) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) &= \mathcal{R}(c_n^*) - \hat{\mathcal{R}}_n(c_n^*) + \hat{\mathcal{R}}_n(c_n^*) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) \\ &\leq \mathcal{R}(c_n^*) - \hat{\mathcal{R}}_n(c_n^*) + \hat{\mathcal{R}}_n(c^*) - \mathcal{R}(c^*) \\ &\leq \sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| + \sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| \\ &\leq 2 \sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)|. \end{aligned}$$

$$\hat{\mathcal{R}}_n(c_n^*) \leq \hat{\mathcal{R}}_n(\tilde{c}) \text{ for all } \tilde{c} \in \mathcal{C}.$$

The result follows. ■

Note that this bound requires uniform convergence, since we minimize \hat{c}_n^* w.r.t. empirical risk. So, the classifier is selected conditional on the data, in which case the weak law of large numbers does not apply.

Theorem 12.8 ([Hoeffding, 1994]). Let X_1, \dots, X_n be independent random variables such that $X_i \in [a_i, b_i]$ and $S_n = \sum_{i=1}^n X_i$, then for any $t > 0$,

$$\begin{aligned}\mathbb{P}(S_n - \mathbb{E}[S_n] \geq t) &\leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \\ \mathbb{P}(S_n - \mathbb{E}[S_n] \leq -t) &\leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).\end{aligned}$$

In the case of normalized sums $\tilde{S}_n = S_n/n$ and $t = n\epsilon$, we have

$$\begin{aligned}\mathbb{P}(\tilde{S}_n - \mathbb{E}[\tilde{S}_n] \geq \epsilon) &\leq \exp\left(-\frac{2n^2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right) \\ &= \exp\left(-\frac{2n\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2/n}\right) \\ &\xrightarrow{n \rightarrow \infty} 0.\end{aligned}$$

We can also bound the absolute deviation,

$$\begin{aligned}\mathbb{P}(|\tilde{S}_n - \mathbb{E}[\tilde{S}_n]| \geq \epsilon) &= \mathbb{P}(\tilde{S}_n - \mathbb{E}[\tilde{S}_n] \geq \epsilon) + \mathbb{P}(\tilde{S}_n - \mathbb{E}[\tilde{S}_n] \leq -\epsilon) \\ &\leq 2 \exp\left(-\frac{2n\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2/n}\right).\end{aligned}$$

The two cases are disjoint.

Theorem 12.9. Assume $|\mathcal{C}| \leq N$, then for all $\epsilon > 0$,

$$\mathbb{P}\left(\sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \epsilon\right) \leq 2N \exp(-2n\epsilon^2).$$

Proof. This is a trivial application of the union bound and Theorem 12.8,

$$\begin{aligned}\mathbb{P}\left(\sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \epsilon\right) &= \mathbb{P}\left(\bigvee_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \epsilon\right) \\ &\leq \sum_{c \in \mathcal{C}} \mathbb{P}(|\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \epsilon) \\ &\leq \sum_{c \in \mathcal{C}} 2 \exp(-2n\epsilon^2) \\ &\leq 2N \exp(-2n\epsilon^2).\end{aligned}$$

Theorem 12.8 and $\hat{\mathcal{R}}_n : \mathcal{C} \rightarrow [0, 1]$.

This concludes the proof. ■

If we define the confidence $\delta \doteq 2N \exp(-2n\epsilon^2)$, then the precision is

$$\epsilon = \sqrt{\frac{\log N - \log(\delta/2)}{2n}}.$$

Thus, we have a $\mathcal{O}(\log |\mathcal{C}|)$ dependence for the precision, which is not bad, unless $|\mathcal{C}| \rightarrow \infty$. In conclusion, we have $\mathcal{R}(c) - \hat{\mathcal{R}}_n(c) \leq \epsilon$ with high probability uniformly for all concepts $c \in \mathcal{C}$,

$$\mathcal{R}(c) \leq \hat{\mathcal{R}}_n(c) + \sqrt{\frac{\log N - \log(\delta/2)}{2n}}.$$

This is the key difference between Theorem 12.6 and Theorem 12.8—Theorem 12.6 only applies to the empirical risk minimizer.

However, the cardinality of \mathcal{C} might be infinite. We can deal with this in two ways,

1. Represent hypotheses by the classifications that they yield;
2. Measure the \mathcal{VC} dimension of \mathcal{C} and select a hypothesis class with finite \mathcal{VC} -dimensionality.

For example, we can apply the first solution to a linear classification problem. There are an infinite number of hyperplane separators, but there are a finite number of ways that the datapoints can be labeled. Instead of treating each hyperplane independently, we bundle all hyperplanes that lead to the same classifications together.

References

- Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- Yoav Freund and Robert E Schapire. A desicion-theoretic generalization of on-line learning and an application to boosting. In *European conference on computational learning theory*, pages 23–37. Springer, 1995.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The collected works of Wassily Hoeffding*, pages 409–426, 1994.
- Jonas Hübötter, Bhavya Sukhija, Lenart Treven, Yarden As, and Andreas Krause. Information-based transductive active learning. *arXiv preprint arXiv:2402.15898*, 2024.
- Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015.
- Pragya Sur and Emmanuel J Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525, 2019.
- Vladimir Vapnik and Alexey Chervonenkis. Theory of pattern recognition, 1974.
- Victor Veitch, Alexander D’Amour, Steve Yadlowsky, and Jacob Eisenstein. Counterfactual invariance to spurious correlations in text classification. *Advances in neural information processing systems*, 34:16196–16208, 2021.
- Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text

compatible image prompt adapter for text-to-image diffusion models.

arXiv preprint arXiv:2308.06721, 2023.