# (Stochastic) Gradient Descent

THE AUSTRALIAN
NATIONAL UNIVERSITY

**Empirical Risk Functional** $R_{\text{emp}}[f] = \dfrac{1}{m} \sum\limits_{i=1}^{m} c(\mathbf{x}_i, y_i, f(\mathbf{x}_i))$

**Idea 1**

Minimize $R_{\text{emp}}[f]$ by performing gradient descent. This leads to

$$f \rightarrow f - \frac{\Lambda}{m} \sum_{i=1}^{m} \partial_f c(\mathbf{x}_i, y_i, f(\mathbf{x}_i))$$

**Problem**

This may be expensive. If the observations are similar, this is very wasteful.

**Idea 2**

Minimize $R_{\text{emp}}[f]$ by performing stochastic gradient descent over the individual terms under the sum.

**Stochastic Gradient** $f \rightarrow f - \Lambda \partial_f c(\mathbf{x}_i, y_i, f(\mathbf{x}_i))$

**Linear Model** $\mathbf{w} \rightarrow \mathbf{w} - \Lambda \mathbf{x}_i c'(\mathbf{x}_i, y_i, f(\mathbf{x}_i))$

# Perceptron Algorithm for Squared Loss

**argument:** Training sample, $\{\mathbf{x}_1, \ldots, \mathbf{x}_m\} \subset \mathcal{X}$, $\{y_1, \ldots, y_m\} \subset \{\pm 1\}$, $\eta$

**returns:** Weight vector $\mathbf{w}$ and threshold $b$.

**function** Perceptron$(X, Y, \eta)$

    **initialize** $\mathbf{w}, b = 0$

    **repeat**

        **for all** $i$ from $i = 1, \ldots, m$

            Compute $f(\mathbf{x}_i) = \left( \left\langle \sum_{l=1}^{i} \alpha_l \Phi(x_l), \Phi(\mathbf{x}_i) \right\rangle + b \right)$

            Update $\mathbf{w}, b$ according to $\mathbf{w}' = \mathbf{w} + \eta \alpha_i \Phi(\mathbf{x}_i)$ and $b' = b + \eta \alpha_i$

            where $\alpha_i = y_i - f(\mathbf{x}_i)$

        **endfor**

    **until** for all $1 \leq i \leq m$ we have $g(\mathbf{x}_i) = y_i$

    **return** $f : \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$

**end**

# Perceptron Algorithm for Huber's Loss

**argument:** Training sample, $\{\mathbf{x}_1, \ldots, \mathbf{x}_m\} \subset \mathcal{X}$, $\{y_1, \ldots, y_m\} \subset \{\pm 1\}$, $\eta$

**returns:** Weight vector $\mathbf{w}$ and threshold $b$.

**function** $\text{Perceptron}(X, Y, \eta)$

    **initialize** $\mathbf{w}, b = 0$

    **repeat**

        **for all** $i$ from $i = 1, \ldots, m$

            Compute $f(\mathbf{x}_i) = \left( \left\langle \sum_{l=1}^{i} \alpha_l \Phi(x_l), \Phi(\mathbf{x}_i) \right\rangle + b \right)$

            Update $\mathbf{w}, b$ according to $\mathbf{w}' = \mathbf{w} + \eta \alpha_i \Phi(\mathbf{x}_i)$ and $b' = b + \eta \alpha_i$

            where $\alpha_i = \begin{cases} \frac{1}{\sigma}(y_i - f(\mathbf{x}_i)) & \text{for } |y_i - f(\mathbf{x}_i)| \leq \sigma \\ \text{sgn}(y_i - f(\mathbf{x}_i)) & \text{otherwise} \end{cases}$

        **endfor**

    **until** for all $1 \leq i \leq m$ we have $g(\mathbf{x}_i) = y_i$

    **return** $f : \mathbf{x} \mapsto \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$

**end**

# Learning Rate

## Classification

For classification, the absolute value of $f$ does not matter. So we need not adjust the learning rate.

## Regression

The absolute value of $f$ is crucial, so we have to get $\eta$ right.

- Large $\eta$: we get quick initial convergence to the target but large fluctuations remain (stochastic gradient can be very noisy).
- Small $\eta$: slow initial convergence to the target but we have a much better quality estimate in the later stages.

## Trick

Make $\eta$ a variable of the time. One can show that $\eta(t) = O(t^{-1})$ is optimal in many cases. This yields quick initial convergence and low fluctuations later.

## Warning

If $f$ is fluctuating, choosing $\eta$ too small will not be useful.

# Maximum Likelihood and Noise Models

**Basic Idea**

We assume that the observations $y_i$ are derived from $f(\mathbf{x}_i)$ by adding noise, i.e. $y_i = f(\mathbf{x}_i) + \xi_i$ where $\xi_i$ is a random variable with density $p(\xi_i)$.

This also means that once we know the type of noise we are dealing with, we may compute conditional densities $p(y|\mathbf{x})$ under the model assumptions.

**Likelihood** $p(Y|f, X) = p((y_1 - f(\mathbf{x}_1)), \ldots, (y_m, f(\mathbf{x}_m)))$

We make the assumption of iid data (to keep the equations simple). This leads to the likelihood

$$\mathcal{L} = \prod_{i=1}^{m} p(y_i - f(\mathbf{x}_i))$$

**Caveat**

The estimates we obtain are only as good as our initial assumptions regarding the type of function expansion and noise. This means that we may not take $p(Y|X)$ at book value.

# Log-Likelihood and Loss Function

## Idea

Log likelihhood and loss function look suspiciously similar, maybe we can find a link .... For simplicity we assume that the that is generated iid.

## Comparison

$$-\mathcal{L}[f] = \sum_{i=1}^{m} \log p(y_i - f(\mathbf{x}_i))$$

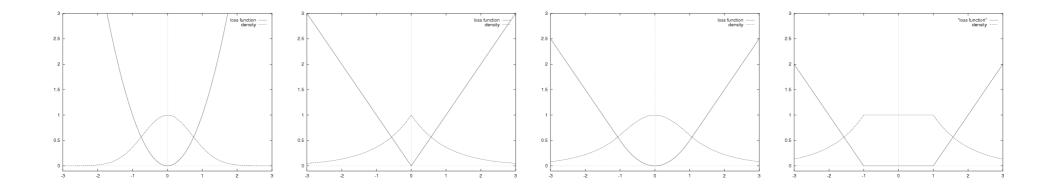$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^{m} c(\mathbf{x}_i, y_i, f(\mathbf{x}_i))$$

## Idea

The two terms differ only by a scaling constant which is irrelevant for minimization purposes. So match up the terms.

$$c(\mathbf{x}, y, f(\mathbf{x})) \equiv -\log p(y_i - f(\mathbf{x}_i))$$

$$p(y_i | f(\mathbf{x}_i) \equiv \exp(-c(\mathbf{x}_i, y_i, f(\mathbf{x}_i))$$

# Density and Loss

|  | loss function $\tilde{c}(\xi)$ | density model $p(\xi)$ |
|---|---|---|
| $\varepsilon$–insensitive | $|\xi|_\varepsilon$ | $\frac{1}{2(1+\varepsilon)} \exp(-|\xi|_\varepsilon)$ |
| Laplacian | $|\xi|$ | $\frac{1}{2} \exp(-|\xi|)$ |
| Gaussian | $\frac{1}{2}\xi^2$ | $\frac{1}{\sqrt{2\pi}} \exp(-\frac{\xi^2}{2})$ |
| Huber's robust loss | $\begin{cases} \frac{1}{2\sigma}(\xi)^2 & \text{if } |\xi| \leq \sigma \\ |\xi| - \frac{\sigma}{2} & \text{otherwise} \end{cases}$ | $\propto \begin{cases} \exp(-\frac{\xi^2}{2\sigma}) & \text{if } |\xi| \leq \sigma \\ \exp(\frac{\sigma}{2} - |\xi|) & \text{otherwise} \end{cases}$ |

THE AUSTRALIAN
NATIONAL UNIVERSITY

**Function Expansion**

We use a linear model (as in the previous lecture) $f_1, \ldots, f_n$ such that

$$f(\mathbf{x}) = \sum_{i=1}^{n} \alpha_i f_i(\mathbf{x})$$

**Additive Noise**

Assume Gaussian noise $\xi$ which corrupts the measurements such that we observe $y$ rather than $f(\mathbf{x})$, i.e. $y = f(\mathbf{x}) + \xi$. We write $\xi \sim \mathcal{N}(0, \sigma)$ in order to state that

$$p(\xi) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}\xi^2}.$$

**Density Model**

From above we know that $p(y|\mathbf{x}, \alpha, \sigma)$ is given by

$$p(y|\mathbf{x}, \alpha, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - f(\mathbf{x}))^2\right)$$

THE AUSTRALIAN
NATIONAL UNIVERSITY

## Likelihood

Under the assumption of iid data, the likelihood of observing $Y = \{y_1, \ldots, y_m\}$, given $X = \{\mathbf{x}_1, \ldots, \mathbf{x}_m\}$ can be found as

$$p(Y|X, \alpha, \sigma) = \prod_{i=1}^{m} p(y_i|\mathbf{x}_i, \alpha, \sigma)$$

## Log Likelihood

$$
\begin{aligned}
\mathcal{L} &= \sum_{i=1}^{m} \log p(y_i|\mathbf{x}_i, \alpha, \sigma) \\
&= \sum_{i=1}^{m} \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_i - f(\mathbf{x}_i))^2\right) \\
&= -\frac{m}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^{m}(y_i - f(\mathbf{x}_i))^2
\end{aligned}
$$

THE AUSTRALIAN
NATIONAL UNIVERSITY

## Optimality Criterion

We need a maximum with respect to the parameters $\alpha, \sigma$. The conditions $\partial_\alpha \mathcal{L} = 0$ and $\partial_\sigma \mathcal{L} = 0$ are necessary for this purpose.

## Optimality in $\alpha$ $\partial_\alpha \mathcal{L} = \partial_\alpha \dfrac{1}{2\sigma^2} \|\mathbf{y} - F\alpha\|^2 = \dfrac{1}{\sigma^2}(F^\top F\alpha - F^\top \mathbf{y}) = 0$

Here we defined (as before) $F_{ij} = f_j(\mathbf{x}_i)$. It leads to the standard least mean squares solution $\alpha = (F^\top F)^{-1} F^\top \mathbf{y}$.

## Optimality in $\sigma$

$$\partial_\sigma \mathcal{L} = \frac{m}{\sigma} - \frac{1}{\sigma^2} \sum_{i=1}^{m} (y_i - f(\mathbf{x}_i))^2 = 0$$

Likewise this leads to $\sigma^2 = \dfrac{1}{m} \displaystyle\sum_{i=1}^{m} (y_i - f(\mathbf{x}_i))^2$ which is *empirical* variance given by the model on the training set.

# When Things go wrong with ML

**No fine-grained prior knowledge**

All functions we optimize over are treated as equally likely.

**Not possible to check assumptions**

- Our ML model works if the assumptions are correct. However, it breaks if they are not all satisfied. And it is hard to test them.

- Difficult to integrate alternative estimates.

- Confidence bounds for estimates.

**High dimensional estimates break**

- Overly confident estimates

- Overfitting

- Likelihood diverges: assume $y_i = f(\mathbf{x}_i)$. In this case we would estimate $\sigma = 0$ as the empirical variance. This in turn leads to $\mathcal{L} \to \infty$.

## Problem

The space of the solutions for $f$ is too large if we admit all possible solutions in, say, the span of $f_1, \ldots, f_n$. Moreover we want to **rank** the solutions.

## Idea

Restrict the possible solutions to the set $\Omega[f] \leq c$ where $\Omega[f]$ is some convex function

$$\Omega[f] = \sum_{i=1}^{n} |\alpha_i| \ (\ell_1 \text{ Regularization})$$

$$\Omega[f] = \frac{1}{2} \sum_{i=1}^{n} \alpha_i^2 \ (\ell_2 \text{ Regularization})$$

$$\Omega[f] = \frac{1}{2} \alpha^\top M \alpha \text{ here } M \text{ is a positive semidefinite matrix}$$

# Regularized Risk Functional

## Problem

Restricting $f$ to the subset $\Omega[f] \leq c$ will solve the problem but the optimization problems are sometimes rather difficult to solve.

## Idea

Trade off the size of $\Omega[f]$ with respect to $R_{\text{emp}}[f]$ and minimize the sum of these two terms.

## Definition

For some $\lambda > 0$, also referred to as the regularization constant, the regularized risk functional is given by

$$R_{\text{reg}}[f] = R_{\text{emp}} + \lambda\Omega[f] = \frac{1}{m}\sum_{i=1}^{m} c(\mathbf{x}_i, y_i, f(\mathbf{x}_i)) + \lambda\Omega[f]$$

This is the central quantity in most learning settings. Note that $R_{\text{reg}}[f]$ is convex, provided $R_{\text{emp}}[f]$ and $\Omega[f]$ are.

**Quadratic Loss** $c(\mathbf{x}, y, f(\mathbf{x})) = \dfrac{1}{2}(y - f(\mathbf{x}))^2$

**Linear Model** $f(\mathbf{x}) = \displaystyle\sum_{i=1}^{n} \alpha_i f_i(\mathbf{x})$

$\ell_2$ **Regularizer** $\Omega[f] = \displaystyle\sum_{i=1}^{n} \alpha_i^2$

**Regularized Risk Functional**

$$R_{\mathrm{reg}}[f] = \frac{1}{m}\sum_{i=1}^{m}\frac{1}{2}(y_i - f(\mathbf{x}_i))^2 + \frac{\lambda}{2}\sum_{i=1}^{n}\alpha_i^2 = \frac{1}{2m}\|\mathbf{y} - F\alpha\|^2 + \frac{\lambda}{2}\|\alpha\|^2$$

**Optimality Conditions**

$\partial_\alpha R_{\mathrm{reg}}[f] = \dfrac{1}{m}(-F^\top \mathbf{y} + F^\top F \alpha) + \lambda\alpha = 0$ and therefore $\alpha = (F^\top F + \lambda m \mathbf{1})^{-1} F^\top \mathbf{y}$
This is the same as when we added $\varepsilon$ to the main diagonal to invert matrices or improve their condition!

# A Practical Example



- Training Set

- Regression for $\lambda = 0.1$

- Regression for $\lambda = 1$

- Regression for $\lambda = 10$