

Mathematics for Machine Learning

(70015)

Lecture 1.1 - Linear Regression Intro

Linear regression aims to provide a solution to the supervised learning problem; we are given a dataset of N examples of inputs and expected outputs, with a goal of predicting the correct output for a new input. Examples of this include image classification (such as digit classification) and translation. A curve fitting problem in 1-dimension has an input space $\in \mathbb{R}$, and an output space $\in \mathbb{R}$.

To tackle this problem mathematically, we need to first describe the curve fitting problem mathematically. As each input is associated with a single output, this is equivalent to a function in mathematics. We are given a dataset of N pairs, of inputs and outputs, where $\mathbf{x}_n \in \mathcal{X}$, which is usually \mathbb{R}^D and $y_n \in \mathcal{Y}$ (usually \mathbb{R} in this case); $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$. The goal is to find a function that maps from the input space to the output space **well**; $f: \mathcal{X} \rightarrow \mathcal{Y}$.

We need to first find candidates for functions that can perform the predictions. Functions need to be parameterised, such that some numbers $\boldsymbol{\theta}$ map to a function. From here, we need to pick the ‘best’ function, thus requiring us to define what good and bad functions are. A good function has the property $f(\mathbf{x}_i, \boldsymbol{\theta}^*) \approx y_i$; the output of the function closely matches the outputs of the training points. This can be defined with a **loss function**, for example;

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i, \boldsymbol{\theta}))^2$$

Therefore, a good function is chosen by minimising the loss; $\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$.

Lecture 1.2 - Scalar Differentiation

We can plot the loss against the parameters for a function. This raises two questions; how to change the parameter to make the loss smaller and how we know if we can’t get a better loss. The derivative is defined as the limit of the difference quotient (as usual);

$$f'(x) = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Several examples of this are as follows;

$f(x)$	$f'(x)$
x^n	nx^{n-1}
$\sin(x)$	$\cos(x)$
$\tanh(x)$	$1 - \tanh^2(x)$
$e^x = \exp(x)$	e^x
$\log(x)$	$\frac{1}{x}$

There are also the following rules which combine the basic functions;

- sum rule describes the derivative of sum of two functions

$$(f(x) + g(x))' = f'(x) + g'(x)$$

- product rule similarly, for multiplication

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$$

- chain rule describes how to differentiate functions that are composed

$$(g \circ f)'(x) = (g(f(x)))' = g'(f(x))f'(x)$$

- quotient rule describes division, special case of the product rule

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$$

This tells us how to change the input in the function; the gradient tells us how much the output changes based on an increase in the input. We first compute the derivative function at a point to find the point's gradient. If the gradient is negative, this tells us the function will decrease if we increase the input (hence increase to minimise), and vice versa; decrease for positive gradients - this is the idea behind gradient descent.

Similarly, we know we are done (at a minimum for the loss function) when there is nothing that can be done to lower the output, hence the gradient must be 0. However, this isn't sufficient to tell us that we have reached a minimum, as a maximum also has a gradient of zero. A minimum has a decreasing function followed by an increasing function; hence the gradient of the gradient (second derivative) is positive.

However, this only gives us a local minima. We should be concerned with getting stuck in a local minima (rather than a global minima) when dealing with non-convex functions. Working through a simple example, with a linear regression problem (aiming to find an optimal a) - the final step takes the second derivative to verify we have a minimum;

$$\begin{aligned} f(x) &= a \cdot x \\ L(a) &= \sum_{n=1}^N (f(x_n) - y_n)^2 \\ \frac{dL}{da} &= \sum_{n=1}^N 2(ax_n - y_n)x_n \\ &= \sum_{n=1}^N 2ax_n^2 - 2x_ny_n \\ &= 0 \\ 2a \sum_n x_n^2 &= \sum_n 2x_ny_n \\ a &= \frac{\sum_n 2x_ny_n}{\sum_n x_n^2} \\ \frac{d^2L}{da^2} &= \sum_{n=1}^N 2x_n^2 \\ &\geq 0 \end{aligned}$$

Lecture 1.3 - Scalar-by-vector Differentiation

The previous example is too simple for real applications - we will need to differentiate by more parameters (vectors). Consider the following example, a polynomial, which has 4 vectors parametising it - each vector now corresponds to a single cubic polynomial;

$$f(x) = \theta_3 x^3 + \theta_2 x^2 + \theta_1 x + \theta_0$$

$$= \boldsymbol{\theta}^\top \boldsymbol{\phi}(x)$$

$$\boldsymbol{\phi}(x) = [x^3 \quad x^2 \quad x \quad 1]^\top$$

Our goal still remains to understand how a function changes with our parameter and to characterise what an optimum is for a function of a vector. Both of these change a multi-dimensional problem into many 1-dimensional problems.

We want to change it into a 1-dimensional question; instead of asking about a change to $\boldsymbol{\theta}$, we ask what happens when we move along a particular line / direction. A directional derivative is how much the function changes, when we move in a direction \mathbf{v} ;

$$\nabla_{\mathbf{v}} L(\boldsymbol{\theta}) = \lim_{h \rightarrow 0} \frac{L(\boldsymbol{\theta} + h\mathbf{v}) - L(\boldsymbol{\theta})}{h}$$

The distance that we move away from the starting point ($\boldsymbol{\theta}$) is determined by **both** the norm / scale of the direction vector \mathbf{v} , as well as h . If we can understand how the function changes based on a change in **any** direction, we can fully characterise differentiation with respect to a vector.

Consider the following example, where we deal with two parameters (also notice the second equality holds as the values in **violet** are equal and cancel each other out);

$$\begin{aligned} \nabla_{\mathbf{v}} L(\boldsymbol{\theta}) &= \lim_{h \rightarrow 0} \frac{L(\theta_1 + hv_1, \theta_2 + hv_2) - L(\theta_1, \theta_2)}{h} \\ &= \lim_{h \rightarrow 0} \underbrace{\frac{L(\theta_1 + hv_1, \theta_2 + hv_2) - L(\theta_1, \theta_2 + hv_2)}{h}}_{\text{only change in first parameter}} + \underbrace{\frac{L(\theta_1, \theta_2 + hv_2) - L(\theta_1, \theta_2)}{h}}_{\text{only change in second parameter}} \\ &= \lim_{h \rightarrow 0} \frac{L(\theta_1 + h', \theta_2 + h' \frac{v_2}{v_1}) - L(\theta_1, \theta_2 + h' \frac{v_2}{v_1})}{\frac{h'}{v_1}} + \frac{L(\theta_1, \theta_2 + h'') - L(\theta_1, \theta_2)}{\frac{h''}{v_2}} \\ &= \frac{\partial L}{\partial \theta_1} v_1 + \frac{\partial L}{\partial \theta_2} v_2 \end{aligned}$$

This means that we can find the gradient in any direction with the partial derivatives, as we chose arbitrary v_1, v_2 . With a partial derivative, we change only one coordinate at a time - see the following for a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$;

$$y = f(\mathbf{x})$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}$$

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{i-1}, \mathbf{x}_i + h, x_{i+1}, \dots, x_N) - f(\mathbf{x})}{h}$$

The Jacobian vector collects all partial derivatives into a row vector;

$$\frac{df}{d\mathbf{x}} = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_N} \right] \in \mathbb{R}^{1 \times N}$$

We now want to know which direction to go in, to decrease the value of the function the most. The directional derivative can be written as the inner product;

$$\nabla_{\mathbf{v}} f(\boldsymbol{\theta}) = \frac{df}{d\boldsymbol{\theta}} \mathbf{v} = \left| \frac{df}{d\boldsymbol{\theta}} \right| |\mathbf{v}| \cos \beta$$

We can maximise this by making $\cos \beta$ as large as possible, since the norms of the two vectors are fixed. If we choose a unit vector \mathbf{v} , we want the largest value possible, hence $\cos \beta = 1$, so the angle between the vectors should be zero (such that $\beta = 0$). As such, we move in the direction of the Jacobian / gradient vector.

We can reuse the intuition from the 1-D case, where moving in either direction doesn't change your value - the directional derivative should be zero in **all directions** (hence the zero vector, $\mathbf{0}$). Additionally, to verify it's a minimum, we want the second directional derivative to also be positive in **all directions**.

Lecture 1.4 - Vector-by-vector Differentiation

Recall the motivating example of linear regression;

$$L(\boldsymbol{\theta}) = \sum_{n=1}^N (y_n - \boldsymbol{\phi}(x_n)^T \boldsymbol{\theta})^2 = \|\mathbf{y} - \boldsymbol{\Phi}(\mathbf{X})\boldsymbol{\theta}\|^2$$

We could either manually take partial derivatives of L , which would be laborious, or consider it as a composition of a vector-to-vector function (the **matrix multiplication**) and a vector-to-scalar function (the norm of the vector squared);

$$f(\mathbf{g}(\boldsymbol{\theta}))$$

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

$$\mathbf{g} : \mathbb{R}^E \rightarrow \mathbb{R}^D$$

There is a multivariate chain rule, for scalars it is as follows (where f is a function of a and b , both of which are functions of t);

$$\frac{df(a(t), b(t))}{dt} = \frac{\partial f}{\partial a} \frac{da}{dt} + \frac{\partial f}{\partial b} \frac{db}{dt}$$

This can be generalised as follows, for $\mathbf{g}(t) \in \mathbb{R}^D$, and both are vectors - allowing us to write it as an inner product;

$$\frac{df(\mathbf{g}(t))}{dt} = \sum_{i=1}^D \frac{\partial f}{\partial g_i} \frac{dg_i}{dt} = \underbrace{\frac{df}{d\mathbf{g}}}_{\text{row}} \cdot \underbrace{\frac{d\mathbf{g}}{dt}}_{\text{col}}$$

This only works as we've defined the differentiation of a function with respect to a vector as a row vector - which we will use for the remainder of the course. Similarly, the second part of the sum is the derivative of a column vector, which remains a column vector. Consider the following example, with $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$;

$$\begin{aligned} f(\mathbf{x}) &= f(x_1, x_2) \\ &= x_1^2 + 2x_2 \end{aligned}$$

$$\begin{aligned} \mathbf{x}(t) &= \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} \\ &= \begin{bmatrix} \sin(t) \\ \cos(t) \end{bmatrix} \end{aligned}$$

$$\frac{df}{d\mathbf{x}} \in \mathbb{R}^{1 \times 2}$$

$$\frac{d\mathbf{x}}{dt} \in \mathbb{R}^2$$

$$\frac{df}{dt} = \frac{df}{d\mathbf{x}} \frac{d\mathbf{x}}{dt}$$

$$= \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix}$$

$$= \begin{bmatrix} 2\sin(t) & 2 \end{bmatrix} \begin{bmatrix} \cos(t) \\ -\sin(t) \end{bmatrix}$$

$$= 2\sin(t)(\cos(t) - 1)$$

A similar rule applies when differentiating with respect to a vector (note previously we only did it with respect to a scalar). Similarly, this just requires collecting all the partial derivatives, and the same chain rule applies for $\mathbf{g}(\mathbf{x}) \in \mathbb{R}^D$;

$$\frac{\partial f(\mathbf{g}(\mathbf{x}))}{\partial x_j} = \sum_{i=1}^D \frac{\partial f}{\partial g_i} \frac{dg_i}{dx_j}$$

Once collected, this becomes matrix multiplication, which can be written in vector form;

$$\frac{df(\mathbf{g}(\mathbf{x}))}{d\mathbf{x}} = \underbrace{\frac{df}{d\mathbf{g}}}_{\text{row}} \cdot \underbrace{\frac{d\mathbf{g}}{d\mathbf{x}}}_{\text{mat}}$$

Note that the matrix is the derivative of a column vector (\mathbf{g}) with respect to an input vector \mathbf{x} . We instead put the elements in each row (similar to how we did it for a vector by scalar) - the elements of \mathbf{g} (i) are along the column, and the dimensions of the derivative (j) are along the row. If f gave a vector as an output, we'd end up with a matrix by matrix multiplication.

Consider the following example, where we have $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$, $\mathbf{y} = \mathbf{f}(\mathbf{x}) \in \mathbb{R}^M$, and $\mathbf{x} \in \mathbb{R}^N$;

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, \dots, x_N) \\ \vdots \\ f_M(x_1, \dots, x_N) \end{bmatrix}$$

The collection of all partial derivatives is called a Jacobian matrix;

$$\begin{bmatrix} \frac{dy_1}{d\mathbf{x}} \\ \vdots \\ \frac{dy_M}{d\mathbf{x}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \dots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

In general, a function $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$ has a gradient that is a matrix $\mathbb{R}^{M \times N}$, or has the number of target dimensions (M) \times the number of input dimensions (N);

$$d\mathbf{f}[m, n] = \frac{\partial f_m}{\partial x_n}$$

Recall that for matrix multiplication, the second dimension of the first matrix must match with the first dimension of the second matrix. A function composition is constrained that the output dimension of \mathbf{h} must be the same as the input dimension of \mathbf{g} to compute $\mathbf{g}(\mathbf{h}(\mathbf{x}))$ when we have $\mathbf{f}(\mathbf{x}) = (\mathbf{g} \circ \mathbf{h})(\mathbf{x})$. This ensures the shapes of the chain rule will **always** work out. For $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$, $\mathbf{g} : \mathbb{R}^L \rightarrow \mathbb{R}^M$, and $\mathbf{h} : \mathbb{R}^N \rightarrow \mathbb{R}^L$, we have the following shapes;

$$\underbrace{\frac{d\mathbf{f}}{d\mathbf{x}}}_{M \times N} = \underbrace{\frac{d\mathbf{g}}{d\mathbf{h}}}_{M \times L} \underbrace{\frac{d\mathbf{h}}{d\mathbf{x}}}_{L \times N}$$

Consider the following example, of a matrix multiplication $\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ where we have $\mathbf{A} \in \mathbb{R}^{M \times N}$ and $\mathbf{x} \in \mathbb{R}^N$;

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} \tag{1}$$

$$= \begin{bmatrix} A_{1,1}x_1 + A_{1,2}x_2 + \dots + A_{1,N}x_N \\ \vdots \\ A_{M,1}x_1 + A_{M,2}x_2 + \dots + A_{M,N}x_N \end{bmatrix} \tag{2}$$

$$f_i(\mathbf{x}) = \sum_{k=1}^N A_{i,k}x_k \tag{3}$$

$$\frac{\partial f_i}{\partial x_j} = \sum_k A_{i,k} \frac{\partial x_k}{\partial x_j} \tag{4}$$

$$= \sum_k A_{i,k} \delta_{kj} \tag{5}$$

$$= A_{i,j} \tag{6}$$

Note the following steps in particular;

- (3) write out each scalar in the vector with index notation by writing the sum explicitly allowing us to take the derivative of an arbitrary element in the output vector with respect to an arbitrary element in the input vector
- (4) take differential operator inside by sum rule
- (5) matrix value is constant, so we end up taking a partial derivative of x_k by x_j - partial derivative only changes x_j and keeps everything else constant, hence it will be zero when $k \neq j$ and 1 when $k = j$ (indicator function, δ)
- (6) end up with $A_{i,j}$ as it's the only case the indicator function is non-zero

Therefore, we get the result that $\frac{df}{dx} = \mathbf{A} \in \mathbb{R}^{M \times N}$.

Consider the following, which is similar to the loss function, where $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, $\mathbf{e}, \mathbf{y} \in \mathbb{R}^M$;

$$\begin{aligned}
L(\mathbf{e}) &= \frac{1}{2} \|\mathbf{e}\|^2 \\
&= \frac{1}{2} \mathbf{e}^\top \mathbf{e} \\
\mathbf{e} &= \mathbf{y} - \mathbf{A}\mathbf{x} \\
\frac{dL}{d\mathbf{x}} &= \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \mathbf{x}} \\
\frac{\partial L}{\partial e_i} &= \frac{\partial}{\partial e_i} \sum_j \frac{1}{2} e_j^2 \\
&= \sum_j \frac{1}{2} 2e_j \frac{\partial e_j}{\partial e_i} \\
&= e_i \\
\frac{\partial L}{\partial \mathbf{e}} &= \mathbf{e}^\top \\
\frac{\partial \mathbf{e}}{\partial \mathbf{x}} &= -\mathbf{A} \\
\frac{\partial L}{\partial \mathbf{x}} &= \mathbf{e}^\top (-\mathbf{A}) \\
&= -(\mathbf{y} - \mathbf{A}\mathbf{x})^\top \mathbf{A}
\end{aligned}$$

Lecture 1.5 - Hessians

We still want to get the second derivative, in order to verify that we have a minima and not a maxima. The chain rule we have in place only works with derivatives of scalars, or column vectors (when differentiating with respect to vectors). We've currently defined the vertical axis to be outputs and the horizontal axis for variables we're differentiating by. This convention no longer holds when we are taking second derivatives, as we need to take the second derivative along a row vector;

$$\nabla_v \underbrace{\left[\frac{df}{d\boldsymbol{\theta}} \right]}_{\text{scalar}} = \underbrace{\frac{d}{d\boldsymbol{\theta}} \left[\frac{df}{d\boldsymbol{\theta}} \right]}_{\text{row vector}} \mathbf{v}$$

Solving this in a way that only involves taking derivatives with respect to scalars;

$$\frac{df}{d\boldsymbol{\theta}} \mathbf{v} = \sum_j \frac{\partial f}{\partial \theta_j} v_j \tag{1}$$

$$\frac{\partial}{\partial \theta_i} \left[\frac{df}{d\boldsymbol{\theta}} \right] \mathbf{v} = \sum_j \frac{\partial}{\partial \theta_i} \frac{\partial f}{\partial \theta_j} v_j \tag{2}$$

$$= \sum_j \frac{\partial^2 f}{\partial \theta_i \partial \theta_j} v_j \quad (3)$$

$$\nabla_v \left[\frac{df}{d\boldsymbol{\theta}} \mathbf{v} \right] = \mathbf{v}^\top \mathbf{H} \mathbf{v} \quad (4)$$

This involves the following steps;

1. expand out in terms of summation components (partial derivatives multiplied by vector component)
2. take partial derivative of the scalar, using sum rule
3. obtain a row vector where we multiply against the second partial derivatives
4. the matrix \mathbf{H} (Hessian) contains all second partial derivatives of the function f

We are at the minimum if \mathbf{H} is positive definite ($\forall \mathbf{v} \mathbf{v}^\top \mathbf{H} \mathbf{v} \geq 0$) - positive eigenvalues.

Lecture 2 - Matrix Derivatives and Backpropagation

Defining a function with respect to a number of parameters, and taking a derivative (by doing many simple problems) can be time consuming.

Consider a loss function;

$$L = \frac{1}{2} \|\mathbf{e}\|^2 \quad \mathbf{e} \in \mathbb{R}^N$$

$$\mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{c} \quad \mathbf{A} \in \mathbb{R}^{N \times D}$$

$$\mathbf{c} = \sin \mathbf{x} \quad \mathbf{x} \in \mathbb{R}^D$$

$$\begin{aligned} \frac{\partial L}{\partial x_j} &= \frac{\partial}{\partial x_j} \frac{1}{2} \sum_{n=1}^N \left(y_n - \sum_i A_{n,i} \sin x_i \right)^2 \\ &= \frac{1}{2} \cdot 2 \cdot \sum_{n=1}^N \left(y_n - \sum_i A_{n,i} \sin x_i \right) \cdot \frac{\partial}{\partial x_j} \left[y_n - \sum_i A_{n,i} \sin x_i \right] \end{aligned}$$

$$= \sum_{n=1}^N \left(y_n - \sum_i A_{n,i} \sin x_i \right) \sum_i \frac{\partial}{\partial x_j} (A_{n,i} \sin x_i)$$

$$= \sum_{n=1}^N \left(y_n - \sum_i A_{n,i} \sin x_i \right) \cdot - \sum_i A_{n,i} \cos x_i \frac{\partial x_i}{\partial x_j}$$

$$= \sum_{n=1}^N \left(y_n - \sum_i A_{n,i} \sin x_i \right) \cdot - \sum_i A_{n,i} \cos x_i \delta_{ij}$$

$$= \sum_{n=1}^N \left(y_n - \sum_i A_{n,i} \sin x_i \right) \cdot - A_{n,j} \cos x_j$$

$$\frac{dL}{d\mathbf{x}} = \underbrace{\frac{dL}{d\mathbf{e}}}_{1 \times N} \cdot \underbrace{\frac{d\mathbf{e}}{d\mathbf{c}}}_{N \times D} \cdot \underbrace{\frac{d\mathbf{c}}{d\mathbf{x}}}_{D \times D}$$

$$\frac{\partial L}{\partial e_j} = \frac{\partial}{\partial e_j} \left[\frac{1}{2} \sum_{n=1}^N e_n^2 \right]$$

$$= \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial e_j} e_n^2$$

$$= \sum_n e_n \delta_{nj}$$

$$\begin{aligned}
&= e_j \\
\frac{\partial e_i}{\partial c_j} &= \frac{\partial}{\partial c_j} \left[y_i - \sum_k A_{i,k} c_k \right] \\
&= - \sum_k A_{i,k} \frac{\partial c_k}{\partial c_j} \\
&= - \sum_k A_{i,k} \delta_{kj} \\
&= -A_{i,j}
\end{aligned}$$

collect j s along row and i s down column

$$\begin{aligned}
\frac{d\mathbf{e}}{d\mathbf{c}} &= -\mathbf{A} \\
\frac{\partial c_i}{\partial x_j} &= \frac{\partial}{\partial x_j} \sin x_i \\
&= \cos x_i \frac{\partial x_i}{\partial x_j} \\
&= (\cos x_i) \delta_{ij}
\end{aligned}$$

note

Note that previously we only had sums with the indicator function; in this case, we only have $\cos x_i$ along the diagonal, and zeroes elsewhere.

Consider the linear regression problem;

$$\begin{aligned}
L(\boldsymbol{\theta}) &= \frac{1}{2} \|\mathbf{y} - \Phi(\mathbf{X})\boldsymbol{\theta}\|^2 \\
\boldsymbol{\theta} &\in \mathbb{R}^3 \\
\Phi(\mathbf{X}) &\in \mathbb{R}^{N \times 3} \\
\mathbf{y} &\in \mathbb{R}^N \\
\phi(x) &= \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} \\
\Phi(\mathbf{X}) &= \begin{bmatrix} \phi^\top(x_1) \\ \phi^\top(x_2) \\ \vdots \end{bmatrix} \\
f(x) &= \phi(x)^\top \boldsymbol{\theta} \\
&= \theta_0 + \theta_1 x + \theta_2 x^2 \\
L(\mathbf{e}) &= \frac{1}{2} \|\mathbf{e}\|^2 \\
\mathbf{e} &= \mathbf{y} - \Phi(\mathbf{X})\boldsymbol{\theta} \\
\frac{dL}{d\boldsymbol{\theta}} &= \frac{dL}{d\mathbf{e}} \cdot \frac{d\mathbf{e}}{d\boldsymbol{\theta}} \\
\frac{dL}{d\mathbf{e}} &= \mathbf{e}^\top \\
\frac{d\mathbf{e}}{d\boldsymbol{\theta}} &= -\Phi(\mathbf{X}) \\
&= (\mathbf{y} - \Phi(\mathbf{X})\boldsymbol{\theta})^\top \cdot -\Phi(\mathbf{X}) \\
&= 0 \\
\boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top \Phi(\mathbf{X}) &= \mathbf{y}^\top \Phi(\mathbf{X})
\end{aligned}$$

$$\theta^\top \underbrace{[\Phi(\mathbf{X})^\top \Phi(\mathbf{X})][\Phi(\mathbf{X})^\top \Phi(\mathbf{X})]^{-1}}_I = \mathbf{y}^\top \Phi(\mathbf{X})[\Phi(\mathbf{X})^\top \Phi(\mathbf{X})]^{-1}$$

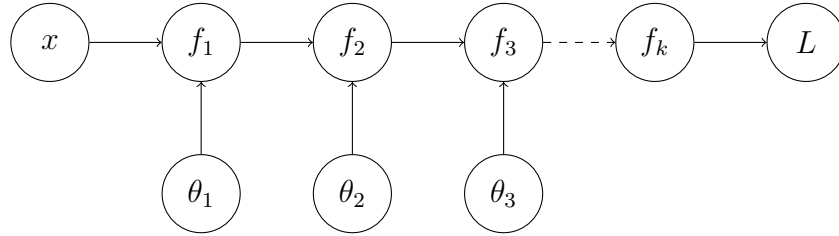
The basis function can also be a learnable feature, with $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^L$ and σ being an activation function;

$$\begin{aligned} \mathbf{x} &\in \mathbb{R}^D \\ \boldsymbol{\theta} &\in \mathbb{R}^L \\ f(\mathbf{x}) &= \phi(\mathbf{x})^\top \boldsymbol{\theta} \\ \phi(x) &= \sigma(\mathbf{A}\mathbf{x} + \mathbf{b}) \end{aligned}$$

This can be repeated deeper by replacing the \mathbf{x} with $\phi'(\mathbf{x})$ (with \mathbf{A}' and \mathbf{b}' similarly), and so on. The idea is that we have a layer that takes in the output of the previous layer as follows (eventually, this stops with f_0 being x , the input);

$$f_l(f_{l-1}) = \sigma(\mathbf{A}_l f_{l-1} + \mathbf{b}_l)$$

Note for brevity, we combine \mathbf{A}_i and \mathbf{b}_i as θ_i ;



If we want to learn in this model, we need to take derivatives with respect to all parameters we've defined (\mathbf{A} and \mathbf{b} in each layer). However, the former is a matrix; hence we need to take the derivative of a function with respect to a matrix, such as (note that this is no longer matrix multiplication, even if they both look like matrices; will be explained later);

$$\frac{d}{d\theta} \mathbf{x}^\top \mathbf{A}(\theta) \mathbf{x} \text{ or } \frac{d}{d\mathbf{A}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$$

When taking the derivative of a vector with respect to a vector, we had two indices which could be nicely represented in a matrix. But in the case of differentiating with respect to a matrix, we have the indices of the vector, as well as a grid of indices in the matrix.

Functions of matrices are no different to functions of vectors; they're still just functions of many different numbers - hence a function of a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ remains a multivariate function.

$$\begin{aligned} f(\mathbf{A}) &= \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 \\ f(A_{1,1}, A_{2,1}, \dots, A_{M,1}, \dots, A_{M,N}) &= \sum_i \left(\sum_{j} A_{i,j} x_j - y_i \right) \end{aligned}$$

The chain rule remains the same; where we compute all the partial derivatives with respect to the elements;

$$\begin{aligned} f(\mathbf{g}) &= \|\mathbf{g}\|^2 \\ \mathbf{g}(\mathbf{A}) &= \mathbf{A}\mathbf{x} - \mathbf{y} \\ \frac{\partial f}{\partial A_{i,j}} &= \sum_k \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial A_{i,j}} \end{aligned}$$

Another example;

$$f(\mathbf{A}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

$$\frac{\partial f}{\partial \theta} = \sum_{j,k} \frac{\partial f}{\partial A_{j,k}} \frac{\partial A_{j,k}}{\partial \theta}$$

The difficulty lies in defining the notation that uses well-defined mathematical operations. The chain rule worked well before as we had a nice multiplication in the form of matrix multiplication; this is no longer the case here (see the note prior). Recall the previous notation;

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta} \text{ and } \frac{df}{d\mathbf{A}} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}}$$

A derivative of a vector with respect to a matrix is a rank 3 tensor (a vector can be seen as a rank 1 tensor, and a matrix as rank 2). This representation with higher-dimensional tensors generalises well. Recall that a function $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$ (M target dimensions and N input dimensions) has the following;

$$\frac{d\mathbf{f}}{d\mathbf{x}} \in \mathbb{R}^{M \times N}, \quad d\mathbf{f}[m, n] = \frac{\partial f_m}{\partial x_n}$$

This generalises when the sizes of the inputs are matrices instead, for example $\mathbf{f} : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$, where the gradient is a tensor;

$$\frac{d\mathbf{f}}{d\mathbf{X}} \in \mathbb{R}^{M \times N}, \quad d\mathbf{f}[p, q, m, n] = \frac{\partial f_{p,q}}{\partial X_{m,n}}$$

Applying the chain rule to a particular matrix valued function - take the example of $f \in \mathbb{R}$ being a scalar function of $\mathbf{A} \in \mathbb{R}^{(P \times Q) \times (N \times M)}$, which is a function of $\boldsymbol{\theta} \in \mathbb{R}^L$;

$$\underbrace{\frac{df}{d\boldsymbol{\theta}}}_{1 \times L} = \underbrace{\frac{df}{d\mathbf{A}}}_{1 \times (N \times M)} \underbrace{\frac{d\mathbf{A}}{d\boldsymbol{\theta}}}_{(N \times M) \times L} \quad (1)$$

$$= \frac{df}{d\text{vec}\mathbf{A}} \cdot \frac{d\text{vec}\mathbf{A}}{d\boldsymbol{\theta}} \quad (2)$$

$$\text{vec} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a \\ c \\ b \\ d \end{bmatrix} \quad (3)$$

- (1) the issue is that we end up with a tensor $\mathbb{R}^{(N \times M) \times L}$
- (2) we're now back to doing differentiating with vectors, leading to matrix multiplication
- (3) note that changing a matrix into a vector gives a column vector of the matrix by columns

Automatic differentiation packages keep track of all these axes in these higher order tensors which contain these derivatives and figure out which axes to sum over by looking at the input shape of the function $f : \mathbb{R}^{N \times M} \rightarrow \mathbb{R}$, consistent with the matrix valued chain rule. The most unambiguous way to deal with this is to write it all as a vector, leading to matrix multiplication.

Example of Multivariate Function

The function is $f : \mathbb{R}^{N \times M} \rightarrow \mathbb{R}^N$, defined as;

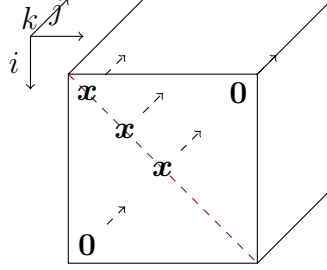
$$\mathbf{f} = \underbrace{\mathbf{A}}_{N \times M} \underbrace{\mathbf{x}}_{M \times 1}$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^{N \times (N \times M)}$$

$$\frac{\partial f_i}{\partial A_{j,k}} = \frac{\partial}{\partial A_{j,k}} \left(\sum_m A_{i,m} x_m \right)$$

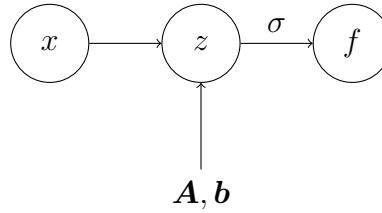
$$\begin{aligned}
&= \sum_m \frac{\partial A_{i,m}}{\partial A_{j,k}} x_m \\
&= \sum_m \delta_{nj} \delta_{mk} x_m \\
&= \delta_{ij} x_k
\end{aligned}$$

This gives the following three-dimensional tensor, taking the convention of putting the outputs (i) along the column. Note that there are zeroes everywhere, other than the **leading diagonal** on i, j where the elements of \mathbf{x} are.



Backpropagation

Recall that a neural network is essentially a stack of function estimators, such that the output of a previous function becomes the input of the next function. The loss function is typically something with respect to an error producing output (observed value y_i compared to some computed value f_i). As such, L is a function of all the parameters that make up the network; $L(\mathbf{A}_1, \mathbf{b}_1, \mathbf{A}_2, \mathbf{b}_2, \dots)$, which can also be written with $\boldsymbol{\theta}$ parameters. Consider the following example of a single-layer network;



In this example, let σ (activation function) be \tanh applied to some linear transformation (which we can denote \mathbf{z}), with $\mathbf{z} \in \mathbb{R}^M$, $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, $\mathbf{b} \in \mathbb{R}^M$. Also note that we have a loss function, L ;

$$\begin{aligned}
\boldsymbol{\theta} &= \{\mathbf{A}, \mathbf{b}\} \\
L(\boldsymbol{\theta}) &= \frac{1}{2} \|\mathbf{e}\|^2 \\
\mathbf{e} &= \mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}) \\
\mathbf{f} &= \tanh(\underbrace{\mathbf{Ax} + \mathbf{b}}_{\mathbf{z}}) \\
\frac{\partial L}{\partial \mathbf{e}} &= \mathbf{e}^\top \\
\frac{\partial \mathbf{e}}{\partial \mathbf{f}} &= -\mathbf{I} \\
\frac{\partial \mathbf{f}}{\partial \mathbf{b}} &= \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}} \\
\frac{\partial \mathbf{f}}{\partial \mathbf{A}} &= \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{A}}
\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathbf{f}}{\partial \mathbf{z}} &= \text{diag}(1 - \tanh^2(\mathbf{z})) \\ \frac{\partial \mathbf{z}}{\partial \mathbf{b}} &= \mathbf{I} \\ \frac{\partial \mathbf{z}}{\partial \mathbf{A}} &= \begin{bmatrix} \mathbf{x}^\top & \cdots & \mathbf{0}^\top & \cdots & \mathbf{0}^\top \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0}^\top & \cdots & \mathbf{x}^\top & \cdots & \mathbf{0}^\top \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0}^\top & \cdots & \mathbf{0}^\top & \cdots & \mathbf{x}^\top \end{bmatrix}\end{aligned}$$

If we blindly applied the rules we know, we'd be able to solve it with many computations of partial derivatives. The idea behind backpropagation is to find an order of computing partial derivatives that maximally reuses what we've computed. Consider the following;

$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K) = \|\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}_K}(\mathbf{f}_{\boldsymbol{\theta}_{K-1}}(\dots \mathbf{f}_{\boldsymbol{\theta}_2}(\mathbf{f}_{\boldsymbol{\theta}_1}(\mathbf{x})) \dots))\|^2$$

This has the following relation (anything that hasn't been previously computed is in **violet**);

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\theta}_K} &= \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \boldsymbol{\theta}_K} \\ \frac{\partial L}{\partial \boldsymbol{\theta}_{K-1}} &= \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \boldsymbol{\theta}_{K-1}} \\ \frac{\partial L}{\partial \boldsymbol{\theta}_{K-2}} &= \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \mathbf{f}_{K-2}} \frac{\partial \mathbf{f}_{K-2}}{\partial \boldsymbol{\theta}_{K-2}} \\ \frac{\partial L}{\partial \boldsymbol{\theta}_{K-3}} &= \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \mathbf{f}_{K-2}} \frac{\partial \mathbf{f}_{K-2}}{\partial \mathbf{f}_{K-3}} \frac{\partial \mathbf{f}_{K-3}}{\partial \boldsymbol{\theta}_{K-3}} \\ \frac{\partial L}{\partial \boldsymbol{\theta}_i} &= \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \dots \frac{\partial \mathbf{f}_{i+1}}{\partial \mathbf{f}_i} \frac{\partial \mathbf{f}_i}{\partial \boldsymbol{\theta}_i}\end{aligned}$$

All the partial gradients that are carried around are vectors, being significantly more computationally efficient than matrices.

Lecture 3 - Probability and Statistics

Overfitting

Recall the curve fitting problem was to find to find a curve that predicts well for unseen inputs, by minimising the loss on the training points. The idea is that it needs to fit the training data well, to fit for unseen data.

$$L(\boldsymbol{\theta}) = \sum_n (f(\mathbf{x}_n; \boldsymbol{\theta}) - y_n)^2$$

Consider how many basis functions we should use (considering polynomial);

$$q \in 0, \dots, Q : \boldsymbol{\phi}(x) = [1 \quad \cdots \quad x^q]^\top$$

For each order we add, we have a more 'wiggly' function that approaches our training data. However, as we get closer to the training data, our predictions further away from our training points start to become more erratic.

$$f_M(\mathbf{x}_n; \boldsymbol{\theta}) = \sum_{m=0}^M x_n^m \theta_m$$

For models $M_1 \leq M_2$, the model M_2 can represent all functions that M_1 can - therefore increasing the order can never worsen the fit on the training data. Closed-form optimisation will find the exact minimum, hence (the training loss will always get smaller when we add basis functions);

$$\min_{\boldsymbol{\theta}} L_{M_2}(\boldsymbol{\theta}) \leq \min_{\boldsymbol{\theta}} L_{M_1}(\boldsymbol{\theta})$$

Typical practice is to split data into a training and a test set. The training set is actually used to find the parameters of a model and the test data shows how well the data predicts on unseen data. We can estimate some loss on the training set, using argmin (which finds the argument that minimises a function);

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} L(\boldsymbol{\theta})$$

$$L_{\text{test}} = \frac{1}{N_{\text{test}}} (f(x_n; \boldsymbol{\theta}^*) - y_n)^2$$

The training error should decrease with higher polynomial degrees, whereas the test error should also decrease until a certain point, when it starts to increase (overfitting).

Probability Theory

We need to formalise this as follows;

- Ω set of all possible outcomes
for a die; $\Omega = \{1, 2, 3, 4, 5, 6\}$
- $E = \{\omega_i \in \Omega\}_{i=1}^I$ event is a set of outcomes
for an event that states the value of the die is greater than 3; $E = \{4, 5, 6\}$
- \mathcal{A} set of all events (a set of sets)
- X, Y random variables (functions on outcomes to the target set τ , simplifies event structure)

$$X : \Omega \rightarrow \tau$$

Consider Ω to be all possible outcomes of a presidential election (can have varying granularity). This could be from the level of what each person voted, and so on. $X(\omega)$ is a simple binary outcome, whether person X is president or not (1 or 0, respectively).

For the scenario of flipping two coins, $\Omega = \{HT, TH, HH, TT\}$. Let $X(\omega)$ be the number of heads;

$$X(\omega) = \begin{cases} 0 & TT \\ 1 & HT, TH \\ 2 & HH \end{cases}$$

A probability measure is a function on sets with the following properties;

- maps from a particular element in the set of all possible events to a real value, larger than 0

$$P : \mathcal{A} \rightarrow \mathbb{R}, P(E) \geq 0$$

- probability of the set of outcomes is equal to 1 (something always happens)

$$P(\Omega) = 1$$

- union of events is the sum when there is no intersection

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_i P(E_i) \text{ when } E_i \cap E_j = \emptyset \forall i, j, i \neq j$$

These axioms describe the reasoning about frequencies of recurring events and the plausibility of beliefs. Our goal is to compute probabilities on the outcomes of random variables (all possible outcomes of Ω

such that $X(\omega)$ is in S - what is the set of all outcomes that would give the random variable values that we care about);

$$P_X(S) = P(X \in S) = P(X^{-1}(S)) = P(\{\omega \in \Omega : X(\omega) \in S\})$$

We can also have multiple random variables as an outcome of the same event. Using the concrete example of the presidential election;

$$X(\omega) = \begin{cases} 1 & \text{person } X \text{ wins} \\ 0 & \text{otherwise} \end{cases}$$

$$Y(\omega) = \begin{cases} 1 & \text{person } X \text{ wins Arkansas} \\ 0 & \text{otherwise} \end{cases}$$

Similarly, we can ask probabilities on the pairs, with the same intuition;

$$P_{X,Y}(S) = P((X,Y) \in S) = P(\{\omega \in \Omega : (X(\omega), Y(\omega)) \in S\})$$

We might want to compute the probability of an event happening for a random variable. For a finite target set, we just need to specify the probability of each of these outcomes ($\frac{1}{6}$ for each die face). Note that for a single event, common notation could be $P(X = x) = p_X(x)$, or even lazier $p(x)$ (with the x implying a random variable X). Similar notation can be used for multiple random variables; $P(X = x, Y = y) = P_{X,Y}(x, y) = P(x, y)$. In the discrete case, probabilities just need to be specified for each pair.

If we have the probabilities of multiple random variables, we can use this to find the probability of just one of the RVs. This is marginal probability. For example, if we have $p_{X,Y}(x, y)$ and we just want $p(x)$;

$$P(X = x) = P((X, Y) \in \{(x, y) : y \in \tau_Y\})$$

$$= \sum_{y \in \tau_Y} p_{X,Y}(x, y)$$

An example of this is from *MacKay*; looking at a certain row of the joint table and normalise the occurrences on the same row to sum to 1. Conditional probability is as follows;

$$P(X = x \mid Y = y) = \frac{P(X = x, Y = y)}{\sum_{x'} P(X = x', Y = y)}$$

$$= \frac{P(X = x, Y = y)}{P(Y = y)}$$

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

The two rules are as follows (and Bayes' rule can be derived);

- sum rule $p(x) = \sum_y p(x, y)$
- product rule $p(x, y) = p(y \mid x)p(x) = p(y)p(x \mid y)$
- Bayes' rule $p(y \mid x) = \frac{p(x \mid y)p(y)}{p(x)}$

Statistical independence states that two events are independent if $P(A \cap B) = P(A)P(B)$ or $P(A \mid B) = P(A)$ - or A doesn't depend on B . Two random variables are independent if $p(X = x, Y = y) = p(x)p(y)$, also written as $X \perp\!\!\!\perp Y$.

However, when it's a continuous random variable, the probability of any particular outcome is zero. This is the case where we have $X : \Omega \rightarrow \mathbb{R}$. However, it's more meaningful to ask $P(a \leq X \leq b)$, which can be specified as follows;

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx = \int_a^b p(x) dx$$

And similarly for multiple random variables;

$$P(a \leq X \leq b, c \leq Y \leq d) = \int_a^b \int_c^d f_{X,Y}(x, y) dy dx$$

The rules are similar;

- sum rule $p(x) = \int_y p(x, y) dy$
- product and Bayes' rule are exactly the same (note we are using the shorthand $f_X(x) = p(x)$ when $x \in \mathbb{R}$)

The general form of a **moment** is as follows;

$$\mathbb{E}_X[g(X)] = \int p(x)g(x) dx$$

For a distribution over a vector, the mean vector works similar to the partial derivatives from before (note \mathbf{a} is just any 'direction');

$$\begin{aligned} p(\mathbf{x}) &= p(x_1, x_2, \dots) \\ \mathbb{E}_X[\mathbf{a}^\top \mathbf{x}] &= \int \mathbf{a}^\top \mathbf{x} p(\mathbf{x}) d\mathbf{x} \\ &= \int \sum_i a_i x_i p(\mathbf{x}) d\mathbf{x} \\ &= \sum_i a_i \int x_i p(\mathbf{x}) d\mathbf{x} \\ &= \sum_i a_i \int x_i p(x_i) p(\{x_j\}_{i \neq j} | x_i) d\mathbf{x} \\ &= \sum_i a_i \int x_i p(x_i) dx_i \underbrace{\int p(\{x_j\}_{i \neq j} | x_i) d\{x\}}_{=1} \\ &= \mathbf{a}^\top \bar{\mathbf{x}} \end{aligned}$$

A similar approach can be taken for variance;

$$\begin{aligned} \mathbb{V}_X[\mathbf{a}^\top \mathbf{x}] &= \int (\mathbf{a}^\top \mathbf{x} - \mathbf{a}^\top \bar{\mathbf{x}}) p(\mathbf{x}) d\mathbf{x} \\ &= \int \left(\sum_i a_i x_i - a_i - \bar{x}_i \right) \left(\sum_j a_j x_j - a_j - \bar{x}_j \right) p(\mathbf{x}) d\mathbf{x} \\ &= \sum_i \sum_j a_i a_j \int (x_i - \bar{x}_i)(x_j - \bar{x}_j) p(x_i, x_j) d\{x_i, x_j\} \\ &= \mathbf{x}^\top \Sigma \mathbf{a} \end{aligned}$$

The intuition for the direction, for variance, is that we look at the variance in the scalar random variables along a given direction.

Lecture 4 - Cross-validation

We can prevent overfitting by choosing a model (polynomial degree) which isn't too flexible. Generalisation loss represents how well our model performs in future predictions. The setting is that we train the model and deploy it for future predictions, making infinite predictions (limit of $N \rightarrow \infty$) in the future. We also incur a penalty for errors (loss function ℓ).

$$L_{\text{test}} = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n; \boldsymbol{\theta}^*), y_n)$$

We assume the pairs (\mathbf{x}_n, y_n) come from the **same distribution** (both training and future data), hence $\mathbf{x}_n, y_n \sim p(\mathbf{x}, y)$. By the law of large numbers (generalisation error) - note that this is a moment of $p(\mathbf{x}, y)$;

$$L_{\text{test}} = \mathbb{E}_{p(\mathbf{x}, y)}[\ell(f(\mathbf{x}; \boldsymbol{\theta}^*), y)]$$

Our goal is to estimate this loss without knowing the true joint distribution. We can get samples from $p(\mathbf{x}, y)$, and construct an estimator for L_{test} from these samples;

$$\hat{L}_{\text{test}} = \frac{1}{N_{\text{test}}} \sum_{n=1}^N \ell(f(\mathbf{x}_n; \boldsymbol{\theta}^*), y_n)$$

We are writing this without knowing what the density function is, just saying that it exists. As the number of testing points goes to infinity, the estimate approaches the true test loss.

Good ML practice splits data into a training set and a test set, both of which are iid samples from this distribution. Only on the training set do we **minimise** the loss. And loss is **measured** on the test set. As long as **no decisions** are based on the test set, \hat{L}_{test} will be an unbiased estimate of the true average loss. This loss tells us what the error would've been on particular model. Unbiased means;

$$\mathbb{E}_{p(\mathbf{x}_1, \mathbf{x}_2, \dots, y_1, y_2, \dots)}[\hat{L}_{\text{test}}(\mathbf{x}_1, \mathbf{x}_2, \dots, y_1, y_2, \dots)] = L_{\text{test}}$$

We can analyse the variance of this estimator. As the test set becomes larger, the variance decreases.

$$\mathbb{V}_{\Pi_n p(\mathbf{x}_n, y_n)} \left[\frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n; \boldsymbol{\theta}^*), y_n) \right] = \frac{1}{N} \mathbb{V}_{p(\mathbf{x}, y)}[\ell(f(\mathbf{x}; \boldsymbol{\theta}^*), y)]$$

If we have a small variance, the error of the estimator is smaller. Chebyshev's inequality states that the probability of a variable being far away (more than k times the standard deviation) from the mean (the true value we want);

$$P(|X - \bar{X}| \geq k\sigma) \geq \frac{1}{k^2}$$

We want to evaluation a collection of models (different orders of polynomials) on the test set and then pick the model with the lowest test loss. However, we may want to ask whether this is still an unbiased estimate of the test loss. Note that $\boldsymbol{\theta}^*$ is now a function of the test set, therefore the estimator is no longer unbiased - therefore we cannot use the test set to select the model if we want an accurate estimate of the test loss;

$$\mathbb{E}_{\Pi_n p(\mathbf{x}_n, y_n)} \left[\frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n; \boldsymbol{\theta}^*(\{\mathbf{x}_i, y_i\}_{i=1}^N)), y_n) \right] \neq \mathbb{E}_{p(\mathbf{x}, y)}[\ell(f(\mathbf{x}; \boldsymbol{\theta}^*), y)]$$

All labelled data is split into training data and test data. The training data is split further to a training set (to find the parameters of the model) and a validation set (to select the model). Once we've select the model based on the validation set, we can use the true test set (kept completely separate) to find an unbiased estimator. However, with a small validation set there would be large variance (possibly choosing wrong model), whereas a large validation set would lead to smaller training sets, hence wrong parameters could be chosen.

Cross-validation splits the data into training and validation sets in multiple ways, and computes the validation performance for each split. The average of these values is the cross-validation loss (no longer unbiased, however small).

1. split data in K different ways
2. for each model;
 - (a) for each split
 - i. find parameters of model
 - ii. compute loss on validation set
 - (b) calculate average validation loss for all splits
3. pick model with lowest cross-validation loss

Lecture 5 - Gradient Descent

An optimisation is can be formulated as an objective function $L : \mathbb{R}^D \rightarrow \mathbb{R}$ (parameter space to scalar) which roughly tells us how well the model performs on the training data, with an unconstrained minimisation (search parameter in the full space);

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} L(\boldsymbol{\theta})$$

Gradient-based optimisation is a class of methods which perform the following steps;

- pick starting point $\boldsymbol{\theta}_0$
- iterative update the parameters to give a sequence $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T$
- choose the update by computing the gradient $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$
- repeat until we reach a stopping criterion (such as computation budget or gradient size)

The gradient descent algorithm is as follows, with a starting point $\boldsymbol{\theta}_0$ and a sequence of step sizes γ_t ;

1. set $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$
2. repeat the first step until the stopping criterion is met

Consider the following model (note the use of matrices), our goal is to find $\boldsymbol{\theta} \in \mathbb{R}^{D \times 1}$, such that $\mathbf{y} \approx \mathbf{X}\boldsymbol{\theta}$;

$$\begin{aligned} \mathcal{D} &= \{\mathbf{X}, \mathbf{y}\} && \text{dataset} \\ \mathbf{X} &= [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_N]^\top && \in \mathbb{R}^{N \times D} \\ \mathbf{y} &= [y_1 \quad \cdots \quad y_N]^\top && \in \mathbb{R}^{N \times 1} \\ \mathbf{X}^\top &= \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ \downarrow & \downarrow & & \downarrow \\ 1 & 1 & & 1 \end{bmatrix} \end{aligned}$$

The model and loss are as follows (assume a mean of zero and some variance) and the loss is the ℓ_2 norm. We are taking the gradient as the transpose of the derivative of L with respect to $\boldsymbol{\theta}$ (a column vector);

$$\begin{aligned} f(\mathbf{x}, \boldsymbol{\theta}) &= \mathbf{x}^\top \boldsymbol{\theta} \\ y &= f(\mathbf{x}, \boldsymbol{\theta}) + \epsilon && \epsilon \sim \mathcal{N}(0, \sigma^2) \\ L(\boldsymbol{\theta}) &= \frac{1}{2\sigma^2} \sum_n (f(\mathbf{x}_n, \boldsymbol{\theta}) - y_n)^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2 \\
\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) &= \frac{dL}{d\boldsymbol{\theta}} \\
&= \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \quad \text{setting } = 0, \text{ to get pseudo-inverse } \Rightarrow \\
\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta}^* &= \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y} \quad \Rightarrow \\
\boldsymbol{\theta}^* &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}
\end{aligned}$$

The gradient descent for linear regression can be written as follows, assuming constant step-sizes ($\gamma_t = \gamma$);

1. define the starting point $\boldsymbol{\theta}_0$

2. update $\boldsymbol{\theta}_{t+1}$

$$\begin{aligned}
\boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t) \\
&= \boldsymbol{\theta}_t - \gamma \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}_t - \mathbf{y}) \\
&= (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\theta}_t + \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y}
\end{aligned}$$

3. repeat until stopping criterion

We will study the gradient descent update equation (final line). This forms a geometric series ($x_{t+1} = ax_t + b$), which can be written as $\alpha^t x_0 + \beta$ (which applies to vectors as well as scalars).

$$\begin{aligned}
\boldsymbol{\theta}_{t+1} + \beta &= \mathbf{A}(\boldsymbol{\theta}_t + \beta) && \Rightarrow \\
\boldsymbol{\theta}_{t+1} &= \mathbf{A}\boldsymbol{\theta}_t + (\mathbf{A} - \mathbf{I})\beta && \Rightarrow \\
\mathbf{A} &= (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}) \\
(\mathbf{A} - \mathbf{I})\beta &= \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} && \Rightarrow \\
-\frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}\beta &= \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y} && \Rightarrow \\
\beta &= -(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\
&= -\boldsymbol{\theta}^*
\end{aligned}$$

This iterative update gives the following result;

$$\boldsymbol{\theta}_t = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^*$$

Note that this tells us the gradient descent converges (as t gets larger) $\boldsymbol{\theta}_t \rightarrow \boldsymbol{\theta}^*$ if;

$$(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \rightarrow \mathbf{0}$$

In an \mathbb{R}^D space, the space can be constructed by a span of basis vectors. For example, $\mathbb{R}^2 = \text{span}(\{\mathbf{e}_1, \mathbf{e}_2\})$ and $\mathbf{e}_i \perp \mathbf{e}_j, \|\mathbf{e}_i\| = 1$. Any vector in this space can be written as a linear combination of the basis vectors. For the standard basis, assume 2-dimensions, $\mathbf{x} = (x_1, x_2)^\top$, then $\mathbf{x}^\top \mathbf{e}_i = x_i$.

A left multiply of \mathbf{Q}^{-1} on \mathbf{x} is known as a change of basis from $\{\mathbf{e}_1, \mathbf{e}_2\}$ to $\{\mathbf{q}_1, \mathbf{q}_2\}$;

$$\begin{aligned}
\mathbf{Q} &= \begin{bmatrix} q_{1,1} & q_{1,2} \\ q_{2,1} & q_{2,2} \end{bmatrix} \\
&= [\mathbf{q}_1 \quad \mathbf{q}_2]
\end{aligned}$$

$$\begin{aligned}
\mathbf{Q}^{-1} &= \begin{bmatrix} q_{1,1} & q_{2,1} \\ q_{1,2} & q_{2,2} \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{q}_1^\top \\ \mathbf{q}_2^\top \end{bmatrix} \\
\mathbf{z} &= \mathbf{Q}^{-1} \mathbf{x} \\
z_1 &= \mathbf{q}_1^\top \mathbf{x} \\
z_2 &= \mathbf{q}_2^\top \mathbf{x}
\end{aligned}$$

If we were to stretch this vector in the \mathbf{Q} basis;

$$\begin{aligned}
\mathbf{z}' &= \mathbf{\Lambda} \mathbf{Q}^{-1} \mathbf{x} \\
z'_1 &= \lambda_1 \mathbf{q}_1^\top \mathbf{x} \\
z'_2 &= \lambda_2 \mathbf{q}_2^\top \mathbf{x} \\
\mathbf{\Lambda} &= \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}
\end{aligned}$$

To return to the canonical basis, after stretching by $\mathbf{\Lambda}$ in the \mathbf{Q} basis, we have $\mathbf{x}' = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1} \mathbf{x}$. Consider the case where we have $\mathbf{x}' = \mathbf{A} \mathbf{x}$, for some $\mathbf{A} \in \mathbb{R}^{D \times D}$. We want to find $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$. The goal is to find scalar λ (eigenvalue) and vector (eigenvector) \mathbf{q} such that $\mathbf{A} \mathbf{q} = \lambda \mathbf{q}$, if solutions exist for this, we have $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1}$ where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_D)$. If \mathbf{A} is symmetric, there will be D pairs of solutions of eigenvectors and eigenvalues, such that the column vectors ($\mathbf{Q} = (\mathbf{q}_1, \dots, \mathbf{q}_D)$) form an orthonormal basis (such that $\mathbf{Q}^{-1} = \mathbf{Q}^\top$) and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$. Additionally, if \mathbf{A} is positive semi-definite, then $\lambda_D \geq 0$. To find this, we first assume $\mathbf{q} \neq \mathbf{0}$. By $\mathbf{A} \mathbf{q} = \lambda \mathbf{q}$, we can get $(\mathbf{A} - \lambda \mathbf{I}) \mathbf{q} = \mathbf{0}$. With the assumption, we find λ such that $\det(\mathbf{A} - \lambda \mathbf{I}) = 0$ - once we have this, it can be plugged in and solved to obtain \mathbf{q} .

The determinant tells us how the volume is scaled by linear transformation. If the determinant is zero ($\det(\mathbf{A} - \lambda \mathbf{I}) = 0$), then some subspace in \mathbb{R}^D is collapsed into a line $\{\mathbf{0}\}$, implying that there is some $\mathbf{q} \neq \mathbf{0}$ such that $(\mathbf{A} - \lambda \mathbf{I}) \mathbf{q} = \mathbf{0}$. If λ is an eigenvalue of \mathbf{A} , then λ^t is an eigenvalue of \mathbf{A}^t . Additionally, $\lambda + \alpha$ is an eigenvalue of $\mathbf{A} + \alpha \mathbf{I}$;

$$\mathbf{A} + \alpha \mathbf{I} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^{-1} + \mathbf{Q} \alpha \mathbf{I} \mathbf{Q}^{-1} = \mathbf{Q} (\mathbf{\Lambda} + \alpha \mathbf{I}) \mathbf{Q}^{-1} \Rightarrow \mathbf{\Lambda} + \alpha \mathbf{I} = \text{diag}(\lambda_1 + \alpha, \dots, \lambda_D + \alpha)$$

By combining these, if λ is an eigenvalue of \mathbf{A} , then it follows that $(\lambda + \alpha)^t$ is an eigenvalue of $(\mathbf{A} + \alpha \mathbf{I})^t$. Note that the corresponding eigenvector in the pair doesn't change.

The Rayleigh quotient is bounded by the eigenvalues of \mathbf{A} (assuming symmetric);

$$\lambda_{\min}(\mathbf{A}) \leq \underbrace{R(\mathbf{A}, \mathbf{x}) = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\|\mathbf{x}\|_2^2} = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}}_{\text{Rayleigh quotient}} \leq \lambda_{\max}(\mathbf{A}) \Rightarrow \lambda_{\min}(\mathbf{A}) \|\mathbf{x}\|_2^2 \leq \mathbf{x}^\top \mathbf{A} \mathbf{x} \leq \lambda_{\max}(\mathbf{A}) \|\mathbf{x}\|_2^2$$

This result comes from the following;

$$\begin{aligned}
\mathbf{x}^\top \mathbf{A} \mathbf{x} &= \mathbf{x}^\top \mathbf{Q} \mathbf{\Lambda} \underbrace{\mathbf{Q}^\top \mathbf{x}}_{\mathbf{z}} \\
&= \mathbf{z}^\top \mathbf{\Lambda} \mathbf{z} \\
\mathbf{x}^\top \mathbf{x} &= \mathbf{x}^\top \underbrace{\mathbf{Q} \mathbf{Q}^\top}_{\mathbf{I}} \mathbf{x} \\
&= \mathbf{z}^\top \mathbf{z} \\
&= \|\mathbf{z}\|_2^2 \\
&= \sum_{d=1}^D z_d^2
\end{aligned}$$

$$\begin{aligned}
\mathbf{z}^\top \mathbf{\Lambda} \mathbf{z} &= \sum_{d=1}^D \lambda_d z_d^2 \\
R(\mathbf{A}, \mathbf{x}) &= \frac{\mathbf{z}^\top \mathbf{\Lambda} \mathbf{z}}{\|\mathbf{z}\|_2^2} \\
\mathbf{z} &= \begin{bmatrix} z_1 \\ \vdots \\ z_D \end{bmatrix} \\
&= \sum_{d=1}^D \frac{z_d^2}{\|\mathbf{z}\|_2^2} \lambda_d \quad \text{weighted sum}
\end{aligned}$$

Recall the gradient descent with constant step size. The ℓ_2 distance between $\boldsymbol{\theta}_t$ and $\boldsymbol{\theta}^*$;

$$\begin{aligned}
\boldsymbol{\theta}_t &= (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^* \\
\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2^2 &= \|(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)\|_2^2 \\
&= |(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)^\top \underbrace{(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^{2t}}_A \underbrace{(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)}_x|
\end{aligned}$$

Note that this is in the form of the Rayleigh quotient, hence we can bound the distance;

$$\lambda_{\min}((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2)^t \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_2^2 \leq \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2^2 \leq \lambda_{\max}((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2)^t \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*\|_2^2$$

If $\lambda_{\max}^t \rightarrow 0$ as $t \rightarrow \infty$, then we have convergence (as it's an upper bound) as it's a distance between the current solution and the optimal solution. For the rules, we will use the following definitions;

$$\begin{aligned}
\lambda_{\min} &= \lambda_{\min}((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2) & \geq 0 \\
\lambda_{\max} &= \lambda_{\max}((\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2)
\end{aligned}$$

Therefore we have the following rules;

- | | |
|-----------------------------------------------------|------------------------------------------------|
| 1. $\lambda_{\max} < 1$ | always converge |
| 2. $\lambda_{\min} \geq 1$ | always diverge |
| 3. $\lambda_{\min} < 1$ but $\lambda_{\max} \geq 1$ | convergence depends on $\boldsymbol{\theta}_0$ |

To derive the eigenvalues;

- if λ is an eigenvalue of $\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}$ then λ^2 is an eigenvalue of $(\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2$
- if λ is an eigenvalue of $\mathbf{X}^\top \mathbf{X}$ then $1 - \frac{\gamma\lambda}{\sigma^2}$ is an eigenvalue of $\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X}$

This gives us the following result (when combined) - note that $\mathbf{X}^\top \mathbf{X}$ is positive semi-definite, hence $\lambda \geq 0$;

$$\mathbf{X}^\top \mathbf{X} \mathbf{q} = \lambda \mathbf{q} \Leftrightarrow (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})^2 \mathbf{q} = (1 - \frac{\gamma\lambda}{\sigma^2})^2 \mathbf{q}$$

To ensure convergence, we want $\lambda_{\max} < 1$, hence we have the constraint on the step size;

$$\gamma < \frac{2\sigma^2}{\lambda_{\max}(\mathbf{X}^\top \mathbf{X})}$$

Generally, larger step-sizes shouldn't be used as we are unlikely to be lucky (and end up in a divergent case).

1. choose a step-size γ that's larger than the constraint
2. **randomly** initialise the parameter

$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^* + \sum_{d=1}^D \alpha_d \mathbf{q}_d$$

3. iterative update

$$\boldsymbol{\theta}_t - \boldsymbol{\theta}^* = \sum_{d=1}^D \alpha_d \underbrace{\left(1 - \frac{\gamma \lambda_d}{\sigma^2}\right)^t}_k \mathbf{q}_d$$

If $\gamma < \frac{2\sigma^2}{\lambda_d}$ for a direction \mathbf{q}_d , then $k\mathbf{q}_d \rightarrow \mathbf{0}$. Otherwise, $k\mathbf{q}_d$ diverges. Gradient descent will diverge unless $\alpha_d = 0$ for all the diverging directions.

Whether the choice of γ is robust to the initialisation of $\boldsymbol{\theta}_0$ depends on the condition number of $\mathbf{X}^\top \mathbf{X}$.

In general, the loss function isn't quadratic nor is it convex. However, when we approach the optimum, we can take a local quadratic approximation (Taylor expansion of order 2).

Lecture 6 - More on Gradient Descent

We need to check whether the ℓ_2 -norm $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\|_2^2 \rightarrow 0$ when we have $t \rightarrow \infty$. $\boldsymbol{\theta}_t - \boldsymbol{\theta}^* = \mathbf{A}^t(\boldsymbol{\theta}_t - \boldsymbol{\theta}^*)$ where $\mathbf{A} = (\mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X})$. The convergence of gradient descent depends on $\lambda_{\max}(\mathbf{A}^2)$.

The robustness of the choice of step size depends on the condition number $\kappa(\mathbf{X}^\top \mathbf{X})$, with $\kappa \approx 1$ being well conditioned and $\kappa \gg 1$ being ill conditioned (loss function is stretched across some direction). Possible ways of addressing this are as follows;

- **gradient descent with pre-conditioning**

The condition number (mentioned before) is based on the data, and we can't generally control that. The idea of gradient descent with pre-conditioning is as follows, where we apply a linear transformation on the gradient vector;

1. define $\Delta\boldsymbol{\theta}_t = \mathbf{0}$
2. select a pre-conditioner \mathbf{P}_t
3. set $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \mathbf{P}_t^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$
4. set $t \leftarrow t + 1$
5. repeat 2 - 4 until stopping criterion

Consider the following example, for linear regression (assuming constant step sizes γ and fixed pre-conditioner \mathbf{P});

1. set $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma \frac{1}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}_t - \mathbf{y})$
2. set $t \leftarrow t + 1$
3. repeat until stopping criterion

Deriving the geometric sequence form, we see we need to look at the eigenvalues (with the same rules as before) of $(\mathbf{I} - \gamma \frac{1}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^\top \mathbf{X})^2$;

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \gamma \frac{1}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\theta}_t - \mathbf{y}) \\ &= (\mathbf{I} - \gamma \frac{1}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\theta}_t + \gamma \frac{1}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^\top \mathbf{y} \\ \boldsymbol{\theta}_t &= (\mathbf{I} - \gamma \frac{1}{\sigma^2} \mathbf{P}^{-1} \mathbf{X}^\top \mathbf{X})^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) + \boldsymbol{\theta}^* \end{aligned} \Rightarrow$$

As such, to ensure convergence at any initialisation, γ can be chosen as;

$$\gamma < \frac{2\sigma^2}{\lambda_{\max}(\mathbf{P}^{-1}\mathbf{X}^\top\mathbf{X})}$$

Similarly, our robustness depends on the condition number; $\kappa(\mathbf{P}^{-1}\mathbf{X}^\top\mathbf{X})$. As we want a well-conditioned optimisation, we can choose the pre-conditioner to be proportional; $\mathbf{P} \propto \mathbf{X}^\top\mathbf{X}$ - in practice we want $\kappa \approx 1$ with \mathbf{P}^{-1} being easy to compute (expensive to invert). To construct \mathbf{P}_t ;

- use low-rank / diagonal matrices
- have prior knowledge on $\nabla^2 L(\boldsymbol{\theta}_t)$ (useful), for example in linear regression $\nabla^2 L(\boldsymbol{\theta}_t) \propto \mathbf{X}^\top\mathbf{X}$
- $\nabla^2 L(\boldsymbol{\theta}_t)$ may not be constant (or even available) - approximate with statistics of gradients along the update trajectory; used by many adaptive learning rate methods

• gradient descent with momentum

The intuition is to use some historical update information to update the current run. The momentum term is the difference between the current and next parameters. Note that α can depend on t , but it is typically fixed.

1. define $\Delta\boldsymbol{\theta}_t = \mathbf{0}$ and momentum step-size α
2. set $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t) + \alpha \Delta\boldsymbol{\theta}_t$
3. set $\Delta\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t = \alpha \Delta\boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$
4. set $t \leftarrow t + 1$
5. repeat 2 - 4 until stopping criterion

The key idea is to make the current iteration's updates closer to the previous iteration. This helps to speed up the iterations in flatter regions (where $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$ has a small magnitude). It also helps to alleviate oscillations (where $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$ points in a different direction). For the latter, consider the 1-D case; in step t we move to the right, but in $t + 1$, we move to the left (based on our gradient at that point) - the updates cancel slightly to give us an update of a smaller magnitude. Finally, it also helps to smooth out gradients if it's inexact.

• line search

The linear search algorithm automatically adapts γ_t to ensure improvement;

1. define a starting point $\boldsymbol{\theta}_0$ and set $t \leftarrow 0$
2. compute the gradient $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$
3. search $\gamma_t \in (\gamma_{\min}, \gamma_{\max})$ to minimise $L(\boldsymbol{\theta}_t - \gamma \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t))$ with respect to γ
 - **backtracking line search** - shrinking step sizes, let τ be the maximum budget;
 - (a) set a decreasing schedule, $\gamma_{\max} = \gamma^{(1)} > \gamma^{(2)} > \dots > \gamma^{(\tau)} = \gamma_{\min}$
 - (b) try sizes until it satisfies an acceptance criterion, for example $L(\boldsymbol{\theta}_t) - L(\boldsymbol{\theta}_t - \gamma^{(k)} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t))$ is big enough
 - (c) set $\gamma_t \leftarrow \gamma^{(k)}$
 - **other methods** - for example, checking the Wolfe conditions

The decreasing schedule is generally better as we want to do the large movements quickly (imagine a flat region).

4. set $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma_t \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t)$ and $t \leftarrow t + 1$
5. repeat 2 - 4 until stopping criterion

Note that the search space for γ is constrained to ensure we don't go over budget. This is used by nonlinear Conjugate Gradient or BFGS. It converges quickly if $\dim(\boldsymbol{\theta})$ is small.

Stochastic Gradient Descent

All the methods we've discussed require knowing the exact gradient (knowing the loss function and differentiating it). The objective function (LSE) can be written as follows (for regression);

$$\mathbf{g}_t = \nabla_{\boldsymbol{\theta}_t} L(\boldsymbol{\theta}_t) = \nabla_{\boldsymbol{\theta}_t} \sum_{n=1}^N (f(\mathbf{x}_n; \boldsymbol{\theta}_t) - y_n)^2$$

Generally, we can write the loss function as a loss on each individual point;

$$\mathbf{g}_t = \nabla_{\boldsymbol{\theta}_t} L(\boldsymbol{\theta}_t) = \sum_{n=1}^N \nabla_{\boldsymbol{\theta}_t} L_n(\boldsymbol{\theta}_t)$$

However, in big data applications, N could be extremely large. As such, running a full batch gradient descent can be very expensive, as we need to compute L_n for each point as well as storing intermediate results of every $\nabla_{\boldsymbol{\theta}_t} L(\boldsymbol{\theta}_t)$ - just for a single update step.

In SGD, we compute a stochastic estimator for the sum with a random subset $M < N$ terms - the gradient estimator is now random;

$$\hat{\mathbf{g}}_t = \frac{N}{M} \sum_{m \in \mathcal{S}} \nabla_{\boldsymbol{\theta}_t} L_m(\boldsymbol{\theta}_t)$$

If we continuously select "wrong" points, the gradient can end up pointing in the wrong direction. However, by randomising, we don't continue to keep picking the same point.

In practice, this is done by uniform (each point has the same probability of being chosen) sampling with replacement (same point can be chosen again, not taken out of distribution).

$$P(\mathcal{S} = \{m_1, m_2, \dots, m_M\}) = \frac{1}{N} \cdot \frac{1}{N} \cdot \dots = N^{-M}$$

This estimator is unbiased (first condition for convergence);

$$\mathbb{E}_{\mathcal{S}}[\hat{\mathbf{g}}_t] = \mathbf{g}_t \Leftrightarrow \mathbb{E}_{\mathcal{S}} \left[\frac{N}{M} \sum_{m \in \mathcal{S}} \nabla_{\boldsymbol{\theta}_t} L_m(\boldsymbol{\theta}_t) \right] = \sum_{n=1}^N \nabla_{\boldsymbol{\theta}_t} L_n(\boldsymbol{\theta}_t)$$

We also need to choose the step sizes carefully (Robbins-Monro condition), then SGD will converge;

$$\sum_{t=1}^{\infty} \gamma_t = \infty$$
$$\sum_{t=1}^{\infty} \gamma_t^2 < \infty$$

This is a decreasing step size, but not decreasing too quickly - if it's constant, then the solution will swing.

We need to ensure that M is chosen such that it's small enough to be computed quickly, but not too small. The amount of error we will get depends on the error of the stochastic gradient (inversely scaled);

$$\mathbb{V}[\hat{\mathbf{g}}_t] = \mathbb{V}_{\mathcal{S}} \left[\frac{N}{M} \sum_{m \in \mathcal{S}} \nabla_{\boldsymbol{\theta}_t} L_m(\boldsymbol{\theta}_t) \right] = \frac{N^2}{M} \mathbb{V}_{m \sim \text{Uniform}(1, \dots, N)} [\nabla_{\boldsymbol{\theta}_t} L_m(\boldsymbol{\theta}_t)]$$

Lecture 7 - Ridge Regression and PCA

Ridge Regression

Overfitting can be avoided by choosing a model with the right complexity, with the validation data (same distribution with the test data, allowing for generalisation). Another technique to avoid overfitting is regularisation. Our model is as follows, with the goal of finding $\theta \in \mathbb{R}^{D \times 1}$, we assume the noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$ has a known σ^2 ;

$$\begin{aligned}\mathcal{D} &= \{\mathbf{X}, \mathbf{y}\} \\ \mathbf{X} &= [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_N]^\top && \in \mathbb{R}^{N \times D} \\ \mathbf{y} &= [y_1 \quad \cdots \quad y_N]^\top && \in \mathbb{R}^{N \times 1} \\ \mathbf{y} &\approx \mathbf{X}\theta \\ f(\mathbf{x}, \theta) &= \mathbf{x}^\top \theta \\ y &= f(\mathbf{x}; \theta) + \epsilon \\ L(\theta) &= \frac{1}{2\sigma^2} \sum_n (f(\mathbf{x}_n; \theta) - y_n)^2 \\ &= \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|_2^2 && \text{matrix form}\end{aligned}$$

Ridge regression adds an additional **regularisation term** to the loss function, controlled by the hyperparameter λ ;

$$L(\theta) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|_2^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

The optimal solution for θ is θ_R^* , and can be obtained by setting $\nabla_\theta L(\theta) = 0$;

$$\begin{aligned}\theta_R^* &= \operatorname{argmin}_{\theta \in \Theta} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\theta\|_2^2 + \frac{\lambda}{2} \|\theta\|_2^2 \\ \nabla_\theta L(\theta) &= 0 && \Rightarrow \\ \frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X}\theta - \mathbf{y}) + \lambda \theta &= 0 && \Rightarrow \\ \left(\lambda \mathbf{I} + \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} \right) \theta^* &= \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{y} && \Rightarrow \\ \theta_R^* &= (\sigma^2 \lambda \mathbf{I} + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

The optimal solution remains similar to the linear regression case; however the matrix that needs to be inverted is no longer simply $\mathbf{X}^\top \mathbf{X}$. As such, note that if $\lambda = 0$ then $\theta_R^* = \theta^*$. Similarly, gradient descent can be used to solve this, assuming constant step sizes $\gamma_t = \gamma$;

1. define starting point θ_0 , set $t = 0$
2. update $t \leftarrow t + 1$, and update θ_{t+1}

$$\begin{aligned}\theta_{t+1} &= \theta_t - \gamma_t \nabla_\theta L(\theta_t) \\ &= \theta_t - \gamma \left(\frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X}\theta_t - \mathbf{y}) + \lambda \theta_t \right) \\ &= \left((1 - \gamma\lambda) \mathbf{I} - \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{X} \right) \theta_t + \frac{\gamma}{\sigma^2} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

3. repeat until stopping criterion

Non-linear regression uses a non-linear feature mapping $\phi(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^p$ (typically $p > D$). With the non-linear ridge regression model, fitting ridge regression becomes the following;

$$f(\mathbf{x}, \theta) = \phi(\mathbf{x})^\top \theta$$

$$\begin{aligned}
y &= f(\mathbf{x}, \boldsymbol{\theta}) + \epsilon \\
\boldsymbol{\Phi} &= [\phi(\mathbf{x}_1) \quad \cdots \quad \phi(\mathbf{x}_N)]^\top \in \mathbb{R}^{N \times p} \\
\boldsymbol{\theta}_R^* &= \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2 \\
&= (\sigma^2 \lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{y}
\end{aligned}$$

Consider regression with polynomial functions - the parameters we want to fit are the coefficients for each of the polynomial features;

$$f(x, \boldsymbol{\theta}) = \sum_{i=1}^p \theta_i x^{i-1}$$

We can obtain a number of models with different p values that all have a **training loss** close to 0 - however, they may have different **test loss**. Previously, we selected a model based on a validation dataset.

The ℓ_2 regulariser pushes the elements of $\boldsymbol{\theta}$ towards zero. Ridge regression helps balance between data fit and model simplicity; if for a given i , $\theta_i = 0$, then the feature x^{i-1} isn't in use.

$$R(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2 = \sum_{i=1}^p \theta_i^2$$

Ridge regression gives an estimator of $\boldsymbol{\theta}$ with a smaller variance; giving a smaller expected test error than with just linear regression.

We assume there is no model mismatch, hence $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ is generated by $y_n = f(\mathbf{x}_n; \boldsymbol{\theta}_0) + \epsilon_n$. This can also be written as $\mathcal{D} \sim p_{\text{data}}^N$. Note that $\boldsymbol{\theta}_0$ denotes the (unknown) ground truth parameters for the generation of the dataset. We also assume that the test data is also generated from the same process.

Similar to before, our goal is to estimate $\boldsymbol{\theta}$ using data \mathcal{D} , such that $\boldsymbol{\theta}^* \approx \boldsymbol{\theta}_0$ - the solution is written as $\boldsymbol{\theta}^*(\mathcal{D})$, when the dataset is given.

The expected ℓ_2 error for estimate is as follows (ideally we want an unbiased estimator with a low variance);

$$\begin{aligned}
\text{error}(\boldsymbol{\theta}^*) &= \mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [\|\boldsymbol{\theta}^*(\mathcal{D}) - \boldsymbol{\theta}_0\|_2^2] \\
&= \|\boldsymbol{\theta}^*(\mathcal{D}) - \mathbb{E}[\boldsymbol{\theta}^*(\mathcal{D})] + \mathbb{E}[\boldsymbol{\theta}^*(\mathcal{D})] - \boldsymbol{\theta}_0\|_2^2 \\
&= \underbrace{\|\boldsymbol{\theta}^*(\mathcal{D}) - \mathbb{E}[\boldsymbol{\theta}^*(\mathcal{D})]\|_2^2}_{\text{variance}} + \underbrace{\|\mathbb{E}[\boldsymbol{\theta}^*(\mathcal{D})] - \boldsymbol{\theta}_0\|_2^2}_{\text{bias}^2} + \underbrace{\text{cross-term}}_{=0} \\
&= \underbrace{\text{tr}[\mathbb{V}_{\mathcal{D} \sim p_{\text{data}}^N}[\boldsymbol{\theta}^*(\mathcal{D})]]}_{\text{variance}} + \underbrace{\|\mathbb{E}[\boldsymbol{\theta}^*(\mathcal{D})] - \boldsymbol{\theta}_0\|_2^2}_{\text{bias}^2}
\end{aligned}$$

We are not only looking at the estimation error for one particular regression problem (one experiment). To compute the expectation, the experiment is computed multiple times, with each experiment sampling a new set of data from the previously mentioned generation process. The expected error can be written in matrix form;

$$\begin{aligned}
\text{Error}(\boldsymbol{\theta}^*) &= \mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [(\boldsymbol{\theta}^*(\mathcal{D}) - \boldsymbol{\theta}_0)(\boldsymbol{\theta}^*(\mathcal{D}) - \boldsymbol{\theta}_0)^\top] \\
&= \mathbf{b}(\boldsymbol{\theta}^*) \mathbf{b}(\boldsymbol{\theta}^*)^\top \mathbf{V}(\boldsymbol{\theta}^*) \\
\mathbf{b}(\boldsymbol{\theta}^*) &= \mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [\boldsymbol{\theta}^*(\mathcal{D})] - \boldsymbol{\theta}_0 && \text{bias} \\
\mathbf{V}(\boldsymbol{\theta}^*) &= \mathbb{V}_{\mathcal{D} \sim p_{\text{data}}^N} [\boldsymbol{\theta}^*(\mathcal{D})] && \text{variance}
\end{aligned}$$

The expected prediction error is as follows;

$$y_{\text{test}} = \phi(\mathbf{x}_{\text{test}})^\top \boldsymbol{\theta}_0 + \epsilon$$

$$\begin{aligned}
f(\mathbf{x}_{\text{test}}, \boldsymbol{\theta}^*(\mathcal{D})) &= \phi(\mathbf{x}_{\text{test}})^\top \boldsymbol{\theta}^*(\mathcal{D}) \\
\|y_{\text{test}} - f(\mathbf{x}_{\text{test}}, \boldsymbol{\theta}^*(\mathcal{D}))\|_2^2 &= \|\phi(\mathbf{x}_{\text{test}})^\top \boldsymbol{\theta}_0 - \phi(\mathbf{x}_{\text{test}})^\top \boldsymbol{\theta}^*(\mathcal{D}) + \epsilon\|_2^2 \\
&= \|\phi(\mathbf{x}_{\text{test}})^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*(\mathcal{D}))\|_2^2 + \|\epsilon\|_2^2 + \text{cross-term} \\
\mathbb{E}[\|\epsilon\|_2^2] &= \sigma^2 \\
\mathbb{E}[\text{cross-term}] &= 0 \\
\|\phi(\mathbf{x}_{\text{test}})^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*(\mathcal{D}))\|_2^2 &= \phi(\mathbf{x}_{\text{test}})^\top (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*(\mathcal{D})) (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*(\mathcal{D}))^\top \phi(\mathbf{x}_{\text{test}}) \\
\mathbb{E}[(\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*(\mathcal{D})) (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*(\mathcal{D}))^\top] &= \text{Error}(\boldsymbol{\theta}^*) \\
\text{error}_{\text{pred}}(\boldsymbol{\theta}^*) &= \mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [\mathbb{E}_{(\mathbf{x}_{\text{test}}, y_{\text{test}}) \sim p_{\text{data}}^N} [\|y_{\text{test}} - f(\mathbf{x}_{\text{test}}, \boldsymbol{\theta}^*(\mathcal{D}))\|_2^2]] \\
&= \mathbb{E}_{\mathbf{x}_{\text{test}}} [\phi(\mathbf{x}_{\text{test}})^\top \text{Error}(\boldsymbol{\theta}^*) \phi(\mathbf{x}_{\text{test}})] + \sigma^2
\end{aligned}$$

If we have two estimators $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ (for example, the former is the linear regressor and the latter is the ridge regressor) - if we can show the matrix form of the parameter estimator error of the first estimator is dominated by the matrix form of the latter, then we can show the expected prediction error of the former is smaller;

$$\text{Error}(\boldsymbol{\theta}_1) \preceq \text{Error}(\boldsymbol{\theta}_2) \Rightarrow \text{error}_{\text{pred}}(\boldsymbol{\theta}_1) \leq \text{error}_{\text{pred}}(\boldsymbol{\theta}_2)$$

Note that $\text{Error}(\boldsymbol{\theta}_1) \preceq \text{Error}(\boldsymbol{\theta}_2)$ means $\text{Error}(\boldsymbol{\theta}_2) - \text{Error}(\boldsymbol{\theta}_1)$ is positive semi-definite, hence;

$$\forall \phi [\phi^\top (\text{Error}(\boldsymbol{\theta}_1) - \text{Error}(\boldsymbol{\theta}_2)) \phi \leq 0]$$

If we have a smaller estimation error, then we have a smaller prediction error.

Linear regression gives an unbiased estimator (the expectation of ϵ is zero, then the rest cancels out);

$$\mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [\boldsymbol{\theta}_L^*(\mathcal{D})] = \mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [(\Phi^\top \Phi)^{-1} \Phi^\top (\Phi \boldsymbol{\theta}_0 + \epsilon)] = \boldsymbol{\theta}_0$$

However, ridge regression gives a biased estimator (as long as $\lambda \neq 0$);

$$\mathbb{E}_{\mathcal{D} \sim p_{\text{data}}^N} [\boldsymbol{\theta}_R^*(\mathcal{D})] = (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \Phi \boldsymbol{\theta}_0 \neq \boldsymbol{\theta}_0$$

If we know $\epsilon \sim \mathcal{N}(0, \Sigma)$ and $\epsilon' = \mathbf{A}\epsilon$, then $\epsilon' \sim \mathcal{N}(0, \mathbf{A}\Sigma\mathbf{A}^\top)$ The variance (for ridge regression) is as follows;

$$\begin{aligned}
\mathbb{V}_{\mathcal{D} \sim p_{\text{data}}^N} [\boldsymbol{\theta}_R^*(\mathcal{D})] &= \mathbb{V}_{\mathcal{D} \sim p_{\text{data}}^N} [(\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top (\Phi \boldsymbol{\theta}_0 + \epsilon)] \\
&= \mathbb{V}_{\mathcal{D} \sim p_{\text{data}}^N} [(\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \epsilon] \\
&= \sigma^2 (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \Phi (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi) \\
&= \mathbf{V}(\lambda)
\end{aligned}$$

The variance for the linear regression estimator can be obtained from the above, by setting $\lambda = 0$;

$$\mathbf{V}(0) = \sigma^2 (\Phi^\top \Phi)^{-1}$$

The bias of the ridge regression estimator can be written as the follows (and similarly, set $\lambda = 0$ for the bias of the linear regression estimator);

$$\mathbf{b}(\lambda) = -\sigma^2 \lambda (\sigma^2 \lambda \mathbf{I} + \Phi^\top \Phi)^{-1} \boldsymbol{\theta}_0$$

We can choose some λ in the range $0 \leq \lambda \leq \frac{2}{\|\boldsymbol{\theta}_0\|_2^2}$ such that;

$$\mathbf{b}(\lambda) \mathbf{b}(\lambda)^\top + \mathbf{V}(\lambda) \preceq \mathbf{V}(0) \Rightarrow \text{error}_{\text{pred}}(\boldsymbol{\theta}_R^*) \leq \text{error}_{\text{pred}}(\boldsymbol{\theta}_L^*)$$

This shows that ridge regression can achieve a better bias-variance trade-off in the parameter space. For any positive λ , ridge regression gives an estimator that has a smaller variance. However λ needs to be chosen carefully to ensure the bias isn't too large.

Principal Component Analysis

Regression is a supervised learning problem. In unsupervised learning, there is no output data y , and we want to figure out the underlying structure from data X .

The motivation for dimensionality reduction is that the data we observe daily tend to be high dimensional. For example; natural images (actual images) are a very small subset of all possible RGB images - in the majority of cases, it's just noise. Therefore, instead of a high dimensional vector of pixel data, the actual data lies in a lower dimensional space.

A user's data is typically quite sparse in a recommender system; for example a matrix of users' ratings on items - there may be millions of products, but very few will actually have ratings per user. Although the data may seem to be high dimensional, they only lie in a lower dimension subspace.

The problem can be set up as follows (no y , as it's unsupervised);

$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \quad (1)$$

$$\text{mean}(\mathbf{x}_n) = \mathbf{0} \quad (2)$$

$$\mathbf{x}_n \approx \tilde{\mathbf{x}}_n \quad (3)$$

$$= \sum_{j=1}^M z_{nj} \mathbf{b}_j \quad (4)$$

$$z_{nj} = \mathbf{b}_j^\top \mathbf{x}_n \quad (5)$$

$$\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \quad (6)$$

Note the following;

- (1) $\mathbf{x}_n \in \mathbb{R}^{D \times 1}$
- (2) assume that we've already shifted the data
- (3) an approximate projection in a lower-dimensional space, $M < D$
- (5) note that the coordinates $\{z_{nj}\}$ are projections of the \mathbf{x}_n vector onto a given basis.
- (6) an orthonormal basis (which we want to figure out), $\mathbf{b}_m \in \mathbb{R}^{D \times 1}$

The goal is to minimise the ℓ_2 reconstruction error (consider the full orthonormal basis \mathbf{B}_{full});

$$\begin{aligned} \mathbf{B}_{\text{full}} &= [\mathbf{b}_1 \quad \cdots \quad \mathbf{b}_M \quad \mathbf{b}_{M+1} \quad \cdots \quad \mathbf{b}_D] && \text{used in new basis, dropped} \\ \mathbf{x}_n &= \underbrace{\sum_{j=1}^M z_{nj} \mathbf{b}_j}_{\tilde{\mathbf{x}}_n} + \sum_{j=M+1}^D z_{nj} \mathbf{b}_j \\ \mathbf{x}_n - \tilde{\mathbf{x}}_n &= \sum_{j=M+1}^D z_{nj} \mathbf{b}_j \\ L &= \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|_2^2 \\ &= \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D z_{nj} \mathbf{b}_j \right\|_2^2 \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \|z_{nj} \mathbf{b}_j\|_2^2 \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D z_{nj}^2 \end{aligned}$$

Note that the ℓ_2 norm can be put inside the (penultimate step) as the cross terms cancel out, since it is an orthonormal basis. The final step can be done as z_{nj} is a scalar, and the basis vectors are unit vectors. Plugging in the definition of z_{nj} , we get the following (notice the **variance of $\{\mathbf{x}_n\}$**);

$$L = \sum_{j=M+1}^D \mathbf{b}_j^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_j$$

Therefore, our goal is to find the orthonormal basis \mathbf{B}_{full} to minimise L . Note that $\mathbf{S} \in \mathbb{R}^{D \times D}$ is a symmetric matrix and positive semidefinite - for any $\mathbf{x} \in \mathbb{R}^{D \times 1}$; $\mathbf{x}^\top \mathbf{S} \mathbf{x} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}^\top \mathbf{x}_n)^2 \geq 0$.

$$\begin{aligned} \text{define } \mathbf{S} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \\ L &= \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j \end{aligned}$$

Assume the eigen decomposition of $\mathbf{S} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$ (\mathbf{Q} constructs an orthonormal basis of the \mathbb{R}^D space);

$$\begin{aligned} \mathbf{\Lambda} &= \text{diag}(\lambda_1, \dots, \lambda_D) & \lambda_1 \geq \dots \geq \lambda_D \geq 0 \\ \mathbf{b}_j^\top \mathbf{S} \mathbf{b}_j &= \mathbf{b}_j^\top \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top \mathbf{b}_j \\ &= \boldsymbol{\beta}_j^\top \mathbf{\Lambda} \boldsymbol{\beta}_j & \text{defined} \\ &= \sum_{d=1}^D \lambda_d \beta_{jd}^2 \\ \boldsymbol{\beta}_j &= \mathbf{Q}^\top \mathbf{b}_j \\ &= [\beta_{j1}, \dots, \beta_{jD}] \\ &= [\mathbf{q}_1^\top \mathbf{b}_j, \dots, \mathbf{q}_D^\top \mathbf{b}_j] \end{aligned}$$

This gives us a constrained optimisation problem, where we have an orthonormal basis ($\|\mathbf{b}_j\|_2^2 = 1$ and vectors are perpendicular);

$$\min_{\mathbf{B}_{\text{full}}} L = \sum_{j=M+1}^D \sum_{d=1}^D \lambda_d \beta_{jd}^2$$

An iterative approach is as follows; first solve for \mathbf{b}_D , then solve \mathbf{b}_j for $j = D-1, \dots, M+1$ - the constraint still applies that \mathbf{b}_D is a unit vector;

$$\min_{\mathbf{b}_D} \sum_{d=1}^D \lambda_d \beta_{Dd}^2$$

Notice the following, hence it is a weighted sum of the eigenvalues;

$$\begin{aligned} \beta_{Dd} &= \mathbf{b}_D^\top \mathbf{q}_d \\ \sum_{d=1}^D \beta_{Dd}^2 &= 1 \end{aligned}$$

Therefore, $\mathbf{b}_D = \mathbf{q}_D$, by ordering of the eigenvalues - it's just the eigenvector with the smallest eigenvalue. The iterative solution relies on proof by induction;

1. $\mathbf{b}_D = \mathbf{q}_D$
2. for $j = D-1, \dots, M+1$, assume $\mathbf{b}_i = \mathbf{q}_i$, where $i > j$ (already solved the previous optimisation problems)

$$(a) \mathbf{b}_j \perp \mathbf{b}_i, i > j \Rightarrow \mathbf{b}_j = \sum_{d=1}^j \beta_{jd} \mathbf{q}_d$$

By the perpendicular constraint, if we consider \mathbf{b}_j as a linear combination of the \mathbf{q} vectors, it cannot contain the \mathbf{q}_i vectors. If there is a contribution, the inner product wouldn't be equal to zero.

$$(b) \|\mathbf{b}_j\|_2^2 = 1 \Rightarrow \sum_{d=1}^j \beta_{jd}^2 = 1$$

(c) the minimisation problems is as follows, with respect to β_{jd} ;

$$\min_{\beta_j} \sum_{d=1}^j \lambda_d \beta_{jd}^2$$

Therefore, $\mathbf{b}_j = \mathbf{q}_j$ (hence $\beta_{jj} = 1$, $\beta_{jd} = 0$ for $d \neq j$)

3. proof by induction tells us that $\mathbf{b}_j = \mathbf{q}_j$ for $j = M + 1, \dots, D$

As such, we are projecting \mathbf{x}_n to an orthogonal complement space;

$$\text{span}(\{\mathbf{q}_j\}_{j=M+1}^D)^\perp = \{\mathbf{x} \in \mathbb{R}^{D \times 1} : \mathbf{x}^\top \mathbf{q}_j = 0, j = M + 1, \dots, D\}$$

$$\mathbf{x}_n = \underbrace{\sum_{j=1}^M z_{nj} \mathbf{b}_j}_{\tilde{\mathbf{x}}_n} + \sum_{j=M+1}^D z_{nj} \mathbf{q}_j \quad \mathbf{b}_i \perp \mathbf{q}_j$$

$$\tilde{\mathbf{x}}_n \in \text{span}(\{\mathbf{q}_j\}_{j=M+1}^D)^\perp$$

In practice, PCA is done by doing the following;

1. start with the original dataset
2. centre the dataset by subtracting the mean from each data point
3. divide by the standard deviation to make the data unit free, the data will have variance 1 along each of the axes
4. compute eigenvalues and eigenvectors of the data covariance matrix
5. project data onto the principal subspace
6. undo the standardisation and move the projected data back into original space

Lecture 8 - More on PCA

The perspective for PCA covered in the previous lecture was to minimise the reconstruction loss (by minimising the distance between the approximate projection and the original value). Another perspective is to take the maximum variance. Consider data in 2 dimensions, if the data doesn't vary much in the first dimension, then it can be omitted, and only store the other direction. We essentially want to store the dimensions that vary the most, rather than waste storage storing repeated similar values in the other dimensions.

Using a similar set up to before, with an orthonormal basis $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M]$;

$$\mathbf{z}_n := \mathbf{B}^\top \mathbf{x}_n \Leftrightarrow z_{nj} := \mathbf{b}_j^\top \mathbf{x}_n$$

Solve for \mathbf{b}_1 such that $\mathbb{V}[\mathbf{b}_1^\top \mathbf{x}_n]$ is maximised, with the constraint that $\|\mathbf{b}_1\|_2 = 1$. Recall that \mathbf{x}_n has mean zero;

$$\mathbb{V}[\mathbf{b}_1^\top \mathbf{x}_n] = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2$$

$$\begin{aligned}
&= \mathbf{b}_1^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)}_{=\mathbf{S}=\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top} \mathbf{b}_1 \\
&= \mathbf{b}_1^\top \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top \mathbf{b}_1 \\
&= \sum_{d=1}^D \lambda_d \beta_{1d}^2
\end{aligned}$$

Therefore $\mathbf{b}_1 = \mathbf{q}_1$, the eigenvector with the greatest eigenvalue. To show the norm of the $\beta_1 = 1$, we can do the following;

$$\begin{aligned}
\beta_1 &= \mathbf{Q}^\top \mathbf{b}_1 \\
\|\beta_1\|_2^2 &= \|\mathbf{Q}^\top \mathbf{b}_1\|_2^2 \\
&= \left\| \sum_{d=1}^D (\mathbf{q}_d^\top \mathbf{b}_1) \mathbf{q}_d \right\|_2^2 && \mathbf{Q} \text{ contains the basis (column) vectors} \\
&= \left\| \sum_{d=1}^D \beta_{1d} \mathbf{q}_d \right\|_2^2 \\
&= \sum_d \sum_{d'} \beta_{1d} \beta_{1d'} \mathbf{q}_d^\top \mathbf{q}_{d'} \\
&= \sum_d \beta_{1d}^2 \|\mathbf{q}_d\|_2^2 && \text{different } \mathbf{q} \text{ vectors are orthonormal} \\
&= \sum_d \beta_{1d}^2 && \mathbf{q} \text{ vectors are unit vectors}
\end{aligned}$$

The remaining directions, $\mathbf{b}_2, \dots, \mathbf{b}_M$ can be solved iteratively (by ‘pinning’ the directions we’ve already solved) - for $m = 2, \dots, M$;

1. compute the remainder of the problem;

$$\hat{\mathbf{x}}_n = \mathbf{x}_n - \sum_{j=1}^{m-1} z_{nj} \mathbf{b}_j = \mathbf{x}_n - \sum_{j=1}^{m-1} (\mathbf{b}_j^\top \mathbf{x}_n) \mathbf{b}_j$$

2. maximise $\mathbb{V}[\mathbf{b}_m^\top \hat{\mathbf{x}}_n]$, with the constraints that $\|\mathbf{b}_m\|_2 = 1$ and $\mathbf{b}_m \perp \mathbf{b}_j$ for $j = 1, \dots, m-1$ (perpendicular to the previous solutions), the mean is still $\mathbf{0}$, the variance after the projection is as follows;

$$\begin{aligned}
\mathbb{V}[\mathbf{b}_m^\top \hat{\mathbf{x}}_n] &= \frac{1}{N} \sum_{n=1}^N \left(\mathbf{b}_m^\top \mathbf{x}_n - \sum_{j=1}^{m-1} (\mathbf{b}_j^\top \mathbf{x}_n) \mathbf{b}_j^\top \mathbf{b}_m \right)^2 \\
&= \mathbf{b}_m^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)}_{=\mathbf{S}=\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top} \mathbf{b}_m && \mathbf{b}_j^\top \mathbf{b}_m = 0
\end{aligned}$$

The iterative solution once again relies on proof by induction;

1. $\mathbf{b}_1 = \mathbf{q}_1$
2. for $m = 2, \dots, M$, assume $\mathbf{b}_j = \mathbf{q}_j$, where $j = 1, \dots, m-1$ (already solved the previous optimisation problems)

$$(a) \mathbf{b}_m \perp \mathbf{b}_j, j = 1, \dots, m-1 \Rightarrow \mathbf{b}_m = \sum_{j=m}^D \beta_{mj} \mathbf{q}_j$$

Set $\beta_{mj} = 0$ for $j = 1, \dots, m-1$

$$(b) \|\mathbf{b}_m\|_2^2 = 1 \Rightarrow \sum_{j=m}^D \beta_{mj}^2 = 1$$

(c) solve for the maximum of $\mathbb{V}[\mathbf{b}_m^\top \mathbf{x}_n]$ with respect to β_{mj} ;

$$\mathbb{V}[\mathbf{b}_m^\top \mathbf{x}_n] = \left(\sum_{j=m}^D \beta_{mj} \mathbf{q}_j \right)^\top \mathbf{S} \left(\sum_{j=m}^D \beta_{mj} \mathbf{q}_j \right) = \sum_{d=m}^D \lambda_d \beta_{md}^2$$

Therefore, $\mathbf{b}_m = \mathbf{q}_m$ (hence $\beta_{mm} = 1$, $\beta_{mj} = 0$ for $j > m$)

3. proof by induction tells us that $\mathbf{b}_m = \mathbf{q}_m$ for $m = 1, \dots, M$

As such, we are projecting \mathbf{x}_n to a subspace;

$$\begin{aligned} \text{span}(\{\mathbf{q}_m\}_{m=1}^M) &= \text{span}(\{\mathbf{q}_j\}_{j=M+1}^D)^\perp \\ \mathbf{x}_n &= \underbrace{\sum_{j=1}^M z_{nj} \mathbf{q}_j}_{\tilde{\mathbf{x}}_n} + \sum_{j=M+1}^D z_{nj} \mathbf{b}_j \quad \mathbf{b}_i \perp \mathbf{q}_j \\ \tilde{\mathbf{x}}_n &\in \text{span}(\{\mathbf{q}_m\}_{m=1}^M) \end{aligned}$$

In both views of PCA, $\mathbf{b}_i \perp \mathbf{b}_j$ and $\mathbf{b}_i \perp \mathbf{q}_j$ - neither view solves for the full basis.

- **minimum reconstruction error**

$$\mathbf{B}_{\text{full}}^* = \{\mathbf{b}_1, \dots, \mathbf{b}_M, \mathbf{q}_{M+1}, \dots, \mathbf{q}_D\}$$

solve for basis vectors of the remainder information; drop last $D - M$ eigenvectors

- **maximum variance**

$$\mathbf{B}_{\text{full}}^* = \{\mathbf{q}_1, \dots, \mathbf{q}_M, \mathbf{b}_{M+1}, \dots, \mathbf{b}_D\}$$

maintain the information in the span of the first M eigenvectors

In practice, we tend to use $\mathbf{B}_{\text{full}}^* = \mathbf{Q}$.

Singular Value Decomposition

PCA frequently takes the eigen decomposition of the variance;

$$\begin{aligned} \mathbf{S} &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \\ &= \frac{1}{N} \mathbf{X} \mathbf{X}^\top \\ \mathbf{X} &= [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N} \end{aligned}$$

The eigendecomposition is $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$;

$$\begin{aligned} \mathbf{U} &= [\mathbf{u}_1, \dots, \mathbf{u}_D] \in \mathbb{R}^{D \times D} && \text{orthonormal basis of } \mathbb{R}^D \\ \mathbf{V} &= [\mathbf{v}_1, \dots, \mathbf{v}_N] \in \mathbb{R}^{N \times N} && \text{orthonormal basis of } \mathbb{R}^N \\ \mathbf{\Sigma} &= \text{diag}(\sigma_1, \dots, \sigma_r) && \text{contains diagonal block with singular values of } \mathbf{X}, \sigma_i > 0 \\ r &= \text{rank}(\mathbf{X}) \end{aligned}$$

$$\leq \min(D, N)$$

Assume a linear mapping $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$, where $\mathbf{A} \in \mathbb{R}^{N \times D}$. For example, a mapping from $D = 2$ to $N = 3$. There exists an orthonormal basis of \mathbb{R}^D , $\{\mathbf{v}_d\}_{d=1}^D$, which forms an orthogonal basis $\langle \mathbf{A}\mathbf{v}_i, \mathbf{A}\mathbf{v}_j \rangle = 0$, for $i \neq j$.

By rescaling $\{\mathbf{A}\mathbf{v}_i\}$ vectors, and adding more orthogonal vectors, another orthonormal basis of \mathbb{R}^N can be constructed; $\{\mathbf{u}_i\}_{i=1}^N$. Therefore $\mathbf{A}\mathbf{v}_i = \sigma_i \mathbf{u}_i$. Note that $\mathbf{A}\mathbf{x}$ can be constructed as linear combination of $\{\mathbf{A}\mathbf{v}_i\}$ vectors. Computing $\mathbf{A}\mathbf{x}$ with $\mathbf{x} \in \mathbb{R}^D$ represented in the standard basis;

- $\mathbf{V}^\top \mathbf{x}$ computes the coordinates of \mathbf{x} in the $\{\mathbf{v}_d\}_{d=1}^D$ basis
- $\Sigma \mathbf{V}^\top \mathbf{x}$ stretches / shrinks the vector
- $\mathbf{U} \Sigma \mathbf{V}^\top \mathbf{x}$ maps to the output space \mathbb{R}^N whose basis is in $\{\mathbf{u}_n\}_{n=1}^N$ and computes the coordinates back in the standard ($\{\mathbf{e}_n\}_{n=1}^N$) basis

A decomposition of \mathbf{S} is as follows;

$$\begin{aligned} \mathbf{X} &= \mathbf{U} \Sigma \mathbf{V}^\top \\ \mathbf{S} &= \frac{1}{N} \mathbf{X} \mathbf{X}^\top \\ &= \frac{1}{N} \mathbf{U} \Sigma \underbrace{\mathbf{V}^\top \mathbf{V}}_{=\mathbf{I}_{N \times N}} \Sigma^\top \mathbf{U}^\top \\ &= \mathbf{U} \text{diag}\left(\frac{\sigma_1^2}{N}, \dots, \frac{\sigma_r^2}{N}, \underbrace{0, \dots, 0}_{D-r}\right) \mathbf{U}^\top \end{aligned}$$

This is equivalent to the eigendecomposition of $\mathbf{S} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$. PCA requires the \mathbf{U} basis vectors, with the largest singular values.

Lecture 9 - Maximum Likelihood

Recall the function is an inner product between a feature vector and some parameters;

$$f(x) = \phi(x)^\top \boldsymbol{\theta} = \sum_{m=1}^M \phi_m(x) \theta_m$$

Our previous principle was to find the function by minimising a loss function on training points; however we are likely to see overfitting. The data points we observe are points that follow the trend with some noise (unpredictability), and our model doesn't account for the latter. We want to have a model that accounts for this unpredictability.

Probabilistic models treat the data as the outcome of some random variable; therefore the model of the data is a probability distribution. For linear regression, we have $y_n = f(x_n; \boldsymbol{\theta}) + \epsilon_n$ (follows the trend with some noise), where $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$ (iid to some Gaussian distribution). This implies there is some probability distribution on y_n . We want to consider the probability density;

$$p(y_n \mid f(x_n; \boldsymbol{\theta}))$$

Given some value for the trend function, $f(x_n; \boldsymbol{\theta})$, it follows that y_n is a deterministic function of a random variable - a change of variables problem.

The problem we are looking at (for discrete RVs), if we are given some probability mass function $P_x(x)$, we want to find the probability mass function of Y where $Y = t(X)$, assuming that t is invertible. We know that $P_X(x) = P_Y(t(x))$ as the outcomes will always occur together. Therefore, we can say $P_Y(y) = P_X(t^{-1}(y))$, by simply looking up the value. If we have the probability for each possible outcome, we can find the probability of any event; $P(Y \in \{y_1, y_2\}) = P(y_1) + P(y_2)$

Continuous RVs are slightly more complex. We now have a continuous random variable X , with the same setting (with the added constraint that t is increasing). Recall the probability of a single exact outcome is zero in the continuous case, instead we care about ranges, such as $P(a \leq Y \leq b)$. We know that equivalent events should have the same probability;

$$\begin{aligned}
P(-\infty \leq X \leq x) &= P(-\infty \leq Y \leq t(x)) \\
&= \int_{-\infty}^x p_X(\bar{x}) d\bar{x} \\
&= \int_{-\infty}^{t(x)} p_Y(\bar{y}) d\bar{y} \\
&= \int_{-\infty}^{t^{-1}(x)} p_X(\bar{x}) d\bar{x} \\
&= \int_{-\infty}^y p_Y(\bar{y}) d\bar{y} \\
\frac{d}{dy} \int_{-\infty}^y p_Y(\bar{y}) d\bar{y} &= p_Y(y) \\
&= \frac{d}{dy} \int_{-\infty}^{t^{-1}(x)} p_X(\bar{x}) d\bar{x} \\
&= \frac{d}{dx} \left(\int_{-\infty}^{t^{-1}(x)} p_X(\bar{x}) d\bar{x} \right) \frac{dx}{dy} \\
&= p_X(x) \frac{d}{dy} t^{-1}(y)
\end{aligned}$$

If t is decreasing, then $P(-\infty \leq X \leq x) = P(t(x) \leq y \leq \infty)$ - in general $p_Y(y) = p_X(x) \left| \frac{d}{dy} t^{-1}(y) \right|$

To work out the probability density $p(y_n \mid f(x_n; \boldsymbol{\theta}))$;

$$\begin{aligned}
y_n &= t(\epsilon_n) \\
t(\epsilon_n) &= \epsilon_n + f(x_n; \boldsymbol{\theta}) \\
p_{Y_n}(y) &= p_{\epsilon_n}(y - f(x_n; \boldsymbol{\theta})) \\
&= \mathcal{N}(y - f(x_n; \boldsymbol{\theta}); 0, \sigma^2) \\
\mathcal{N}(y; \mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}
\end{aligned}$$

Therefore, $p(y_n \mid f(x_n; \boldsymbol{\theta})) = \mathcal{N}(y_n; f(x_n; \boldsymbol{\theta}), \sigma^2)$. Note that $p(y_n \mid f(x_n; \boldsymbol{\theta}))$ may also be written as $p(y_n \mid \boldsymbol{\theta}, x_n)$ or $p(y_n \mid \boldsymbol{\theta})$, since all x_n are given and fixed (rather than unknown).

Our model has free parameters $\boldsymbol{\theta}$ and σ^2 ; the former we have found through loss minimisation. We want to find a single principal to find both. A good model gives high probability to the data; the maximum likelihood estimator.

Assume the height of a random person is Gaussian with unknown mean μ and variance σ^2 . To find the MLE for N observed heights.

$$\begin{aligned}
p(h_n \mid \mu, \sigma^2) &= \mathcal{N}(h_n; \mu, \sigma^2) \\
p(\mathbf{h} \mid \mu, \sigma^2) &= \prod_{n=1}^N p(h_n \mid \mu, \sigma^2) \\
\mu^*, \sigma^{2*} &= \operatorname{argmax}_{\mu, \sigma^2} p(\mathbf{h} \mid \mu, \sigma^2) \\
&= \operatorname{argmax}_{\mu, \sigma^2} \log p(\mathbf{h} \mid \mu, \sigma^2)
\end{aligned}$$

$$\begin{aligned}
&= \operatorname{argmax}_{\mu, \sigma^2} \sum_n \log p(h_n \mid \mu, \sigma^2) \\
&= \operatorname{argmax}_{\mu, \sigma^2} \sum_n \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (h_n - \mu)^2 \right)
\end{aligned}$$

Take derivatives;

$$\begin{aligned}
\frac{d}{d\mu} \left(\sum_n \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (h_n - \mu)^2 \right) \right) &= 0 & \Rightarrow \\
\sum_n -\frac{1}{2\sigma^2} \cdot 2(h_n - \mu) \cdot -1 &= \sum_n \frac{1}{\sigma^2} (h_n - \mu) & \\
&= 0 & \Rightarrow \\
\sum_n h_n - N\mu &= 0 & \Rightarrow \\
\hat{\mu} &= \frac{1}{N} \sum_n h_n \\
\frac{d}{d\sigma^2} \left(\sum_n \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (h_n - \mu)^2 \right) \right) &= 0 & \Rightarrow \\
\sum_n \left(-\frac{1}{2\sigma^2} - \frac{1}{2} (h_n - \hat{\mu})^2 \frac{d}{d\sigma^2} (\sigma^2)^{-1} \right) &= 0 & \Rightarrow \\
-\frac{N}{2\sigma^2} - \frac{1}{2} \sum_n (h_n - \hat{\mu})^2 (\sigma^2)^{-2} \cdot -1 &= 0 & \Rightarrow \\
\frac{1}{N} \sum_n (h_n - \hat{\mu})^2 &= \sigma^2
\end{aligned}$$

To show the mean estimator is unbiased;

$$\hat{\mu} = \frac{1}{N} \sum_n h_n \quad (7)$$

$$\mathbb{E}_{\mathbf{h}}[\hat{\mu}] = \mu \quad \text{goal to show unbiased} \quad (8)$$

$$= \mathbb{E}_{\mathbf{h}} \left[\frac{1}{N} \sum_n h_n \right] \quad (9)$$

$$= \frac{1}{N} \sum_n \mathbb{E}_{h_n} [h_n] \quad (10)$$

$$= \frac{1}{N} \sum_n \mu \quad (11)$$

$$= \mu \quad \text{therefore unbiased} \quad (12)$$

However, the variance estimator isn't - we can do this by having a single example where the expectation $\mathbb{E}[\sigma^2]$ is different from the actual variance of the distribution it came from. For the extreme case, with one data point - the mean is the value of a single data point. As such, the variance estimator would be zero, hence it is biased.

The same idea can be used to train a model;

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{y} \mid \boldsymbol{\theta}, X) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{n=1}^N p(y_n \mid \boldsymbol{\theta}, \mathbf{x}_n)$$

Once again, logs can be taken (doesn't change the optima, but easier to compute and more stable than multiplying in floating point);

$$\begin{aligned}
\boldsymbol{\theta}^* &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log p(\mathbf{y} \mid \boldsymbol{\theta}, X) \\
&= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{n=1}^N \log p(y_n \mid \boldsymbol{\theta}, \mathbf{x}_n) \\
&= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_n \left(-\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y_n - f(\mathbf{x}_n; \boldsymbol{\theta}))^2 \right) \\
&= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_n (y_n - f(\mathbf{x}_n; \boldsymbol{\theta}))^2
\end{aligned}$$

The last line can be done as only the squared term is dependent on $\boldsymbol{\theta}$ - notice that this is the same as minimising loss. The additional capability of our model is that we can also estimate the unpredictability of the model. One assumption we are making is that the uncertainty is the same throughout the entire input range of the model. The variance can be shown to be the average squared deviation from the predicted function value (whereas before, in the height example, it was the deviation from $\hat{\mu}$).

When we overfit, the predicted variance of our model approaches zero. This is bad - when we make a new prediction, the uncertainty helps us avoid making an overconfident prediction. In cross-validation, we train our model on the same training data and then use the validation set to pick the model with the lowest validation error; this may make the error on the training data worse, but better on the validation data. This will give us a larger variance.

Bayesian Inference

Consider the example of a coin tossing game; where if you guess correctly you are given a certain amount, and if you guess incorrectly, you pay a certain amount. Knowledge of the coin's probability is important. Now assume you've seen the game played, and the coin has been flipped to be head three times in a row. We assume the following;

- the coin has a property that determines whether it lands on heads or tails, the probability of which won't change
- no inherent 'skill' in tossing the coin - it depends solely on the coin

We are interested in quantities we don't know about (similar to linear regression, when we didn't know the parameters of the function). Previously, the methods we looked at only attempted to come up with a single best estimate, whereas Bayesian Inference quantifies the uncertainty with a probability distribution. Bayes rule can be used to find the uncertainty after observing data. In the coin example;

$s_i \in \{0, 1\}$	coin toss outcome, 1 for heads
$\mathbf{s} \in \{0, 1\}^N$	sequence of outcomes, a binary vector
$h \in \mathbb{R}$	probability of heads

If we were using MLE, all we would have to do is maximise $p(\mathbf{s} \mid h)$, where \mathbf{s} is data we have and h is the unknown;

$N_H = \text{number of heads in } \mathbf{s}$	
$N_T = \text{number of tails in } \mathbf{s}$	
$N = N_H + N_T$	total number of throws
$p(\mathbf{s} \mid h) = h^{N_H} \cdot (1 - h)^{N_T}$	

The belief of what h is before anything has been observed is the prior distribution.

$$p(h \mid \mathbf{s}) = \frac{p(\mathbf{s}, h)}{p(\mathbf{s})} = \frac{p(\mathbf{s} \mid h)p(h)}{p(\mathbf{s})}$$

Note that $p(h)$ is a probability density on a probability. The posterior distribution is a uniform distribution multiplied by a likelihood function (a linear function). Consider a case where we have 2 heads and then 1 tails, hence $\mathbf{s} = [1, 1, 0]$. The probability distribution is a polynomial, where either extremes are now zero, with a peak at $\frac{2}{3}$. However, when we get more data, the distribution should become sharper at 0.5 (with a fair coin), however, if we kept observing tails for example, we'd likely have a sharper peak at zero.

In short, we keep track of uncertainty about all possible values, and use a probability distribution to quantify this uncertainty. This helps to avoid the overconfidence issues from MLE.