

# Machine Learning for Physicists: Recitation Notes

Clark Miyamoto (cm6627@nyu.edu)

January 22, 2026

## Contents

<b>1</b>	<b>Review of Linear Algebra</b>	<b>3</b>
1.1	Singular Value Decomposition . . . . .	3
1.1.1	Definitions . . . . .	3
1.1.2	Illustration of SVD . . . . .	4
1.1.3	Application, inference of signals . . . . .	4
1.2	Matrix Calculus . . . . .	5
1.2.1	Matrix Inversions . . . . .	6
1.3	Time Complexity . . . . .	7
1.4	References . . . . .	7
<b>2</b>	<b>PyTorch 101</b>	<b>8</b>
2.1	Training Dynamics . . . . .	8
2.2	Whitening . . . . .	8
2.3	Back propagation . . . . .	8
2.4	Weight Initialization . . . . .	8
2.5	PyTorch Tutorial . . . . .	8
<b>3</b>	<b>Review of Probability</b>	<b>9</b>
<b>4</b>	<b>Review of Statistics &amp; Loss Functions</b>	<b>9</b>
4.1	Maximum Likelihood Inference & Mean Squared Error . . . . .	10
4.2	Cross Entropy & Another MLE . . . . .	10
4.3	L2 Regularization . . . . .	10
4.4	Minimizing the loss function . . . . .	11
<b>5</b>	<b>Linear Regression</b>	<b>11</b>
5.1	Frequentist, Maximum Likelihood Estimator . . . . .	12
5.2	Bayesian Linear Regression . . . . .	13
<b>6</b>	<b>Double Descent</b>	<b>13</b>
6.1	Soft Inductive Biases . . . . .	13
<b>7</b>	<b>Training Large Models</b>	<b>13</b>
7.1	Transformers . . . . .	13
7.2	$\mu$ P Optimizer . . . . .	13
<b>8</b>	<b>Geometric Deep Learning</b>	<b>13</b>

<b>I</b>	<b>Generative Models</b>	<b>13</b>
<b>9</b>	<b>"Old" Generative Models</b>	<b>13</b>
9.1	Variational Autoencoders . . . . .	13
9.2	Generative Adversarial Networks (GANs) . . . . .	13
9.3	Denoising Diffusion Probablistic Models (DDPMs) . . . . .	13
<b>10</b>	<b>Modern Generative Models</b>	<b>13</b>
10.1	Review of Non-Equilibrium Statistical Mechanics . . . . .	13
10.2	Measure Transport . . . . .	14
10.3	Score Based Diffusion . . . . .	14
10.4	Flow Matching . . . . .	14
10.5	Stochastic Interpolants . . . . .	14

# 1 Review of Linear Algebra

## 1.1 Singular Value Decomposition

Recall the eigen-decomposition of a matrix. Given a symmetric square matrix  $A \in \mathbb{R}^{d \times d}$  with eigenvalues  $\{\lambda_i\}_i$  and eigenvectors  $\{e_i\}_i$ . The matrix could be re-expressed as

$$A = U \Lambda U^T \quad (1.1)$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d) \in \mathbb{R}^{d \times d}$  and  $U \in \mathbb{R}^{d \times d}$  is a matrix whose columns are  $\{e_i\}_i$ .

This decomposition had a lot of nice properties. In particular,  $\Lambda$  is diagonal and  $U$  is orthogonal. This allowed us to do all sorts of stuff easily; for example, matrix power.

What happens if we want to do this on non-symmetric, or even non-square matrices? Well we can use the singular value decomposition (SVD).

### 1.1.1 Definitions

**Definition 1 (Singular Values)** Let  $A \in \mathbb{R}^{m \times n}$ . Now consider  $A^T A \in \mathbb{R}^{n \times n}$ . This is a symmetric matrix so it has positive eigenvalues  $0 \leq \lambda_1 \leq \dots \leq \lambda_n$ . The singular values  $\sigma_i$  for matrix  $A$  are defined as

$$\sigma_i \equiv \sqrt{\lambda_i}, \text{ s.t. } 0 \leq \lambda_1 \leq \dots \leq \lambda_n \quad (1.2)$$

**Fact 1** The number of non-zero singular values of  $A$  correspond to the rank of  $A$ .

*Proof:* Let  $A : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a linear map. Recall by Rank-Nullity theorem  $\text{rank}(A) + \dim \text{Ker}(A) = \dim(\mathbb{R}^d)$ . Recall  $\text{Ker}(A) = \{v : A(v) = 0\}$ , so the dimension of the kernel is the number of zero eigenvalues.

Also notice that  $\text{Ker}(A) = \text{Ker}(A^T A)$ . ( $\implies$ ) Let  $v \in \text{Ker}(A)$ , then  $A^T A v = 0$ , therefore  $v \in \text{Ker}(A^T A)$ . ( $\impliedby$ ) Let  $v \in \text{Ker}(A^T A)$ , then  $A^T A v = 0$ , meaning  $x^T A^T A v = \|A v\|^2 = 0$ , the vector norm is only zero when the vector is zero, therefore  $A v = 0$ , implying  $v \in \text{Ker}(A)$ .

This means the  $\text{rank}(A) = \dim(\mathbb{R}^d) - \dim \text{Ker}(A^T A)$ . The dimension of the kernel of  $A^T A$  is the number of zero singular values of  $A$ . □

**Definition 2 (SVD)**  $A \in \mathbb{R}^{m \times n}$  with singular values  $0 \leq \sigma_1 \leq \dots \leq \sigma_n$ . Let  $r$  denote the rank, or equivalently the number of singular values of  $A$ . The SVD of  $A$  is a decomposition

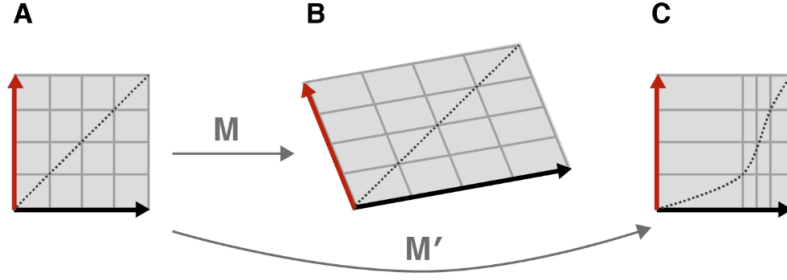
$$A = U \Sigma V^T \quad (1.3)$$

where

- $U \in \mathbb{R}^{m \times m}$  orthogonal matrix
- $V \in \mathbb{R}^{n \times n}$  orthogonal matrix
- $\Sigma \in \mathbb{R}^{m \times n}$  matrix such that  $[\Sigma]_{ii} = \sigma_i$  for  $i \in [1, \dots, r]$  and  $[\Sigma]_{ii} = 0$  for  $i > r$ .

**Theorem 1 (Computing SVD)** Let  $A \in \mathbb{R}^{m \times n}$ . Then  $A$  has a (non-unique) SVD  $A = U \Sigma V^T$ , where

- The columns of  $V$  are orthonormal eigenvectors of  $A^T A$ , where  $A^T A v_i = \sigma_i^2 v_i$ .
- If  $i \leq r$ , s.t.  $\sigma_i \neq 0$ , then the  $i$ 'th column of  $U$  is given by  $\sigma_i^{-1} A v_i$



**Figure 3:** (A) Our original square under a linear transformation  $\mathbf{M}$  (B) and a nonlinear transformation  $\mathbf{M}'$  (C).

Figure 1: Visualization from [2]

### 1.1.2 Illustration of SVD

Recall, by [3Blue1Brown](#), matrix operations transform the coordinate space of some vector. So the only hope we have at visualize the SVD is to use this intuition.

Let  $A \in \mathbb{R}^{2 \times 2}$ . Let  $v_1, v_2 \in \mathbb{R}^2$  be orthonormal vectors, that is  $v_i \cdot v_j = \delta_{ij}$ . Say we know how  $A$  acts on  $v_i$ , that is it rotates them to another orthonormal basis  $u_i$  and rescales them according to  $\sigma_i$ .

$$Av_1 = \sigma_1 u_1 \tag{1.4}$$

$$Av_2 = \sigma_2 u_2 \tag{1.5}$$

We've chosen the notations of these vectors in a very peculiar manner. You can think the SVD as just saying we map vectors in matrix  $V$  to vectors in  $U$  scaled by  $\Sigma$ .

But deriving is believing, so let's show that the matrix  $A$  emits an SVD where everything lines up.

Consider how  $A$  acts on an arbitrary test vector  $x$ .

$$Mx = M(\langle v_1, x \rangle v_1 + \langle v_2, x \rangle v_2) \tag{1.6}$$

$$= Mv_1 \langle v_1, x \rangle + Mv_2 \langle v_2, x \rangle \tag{1.7}$$

$$= \sigma_1 u_1 \langle v_1, x \rangle + \sigma_2 u_2 \langle v_2, x \rangle \tag{1.8}$$

$$= \sigma_1 u_1 v_1^T x + \sigma_2 u_2 v_2^T x \tag{1.9}$$

$$= (\sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T) x \tag{1.10}$$

Since this holds for an arbitrary test vector

$$M = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T = \underbrace{\begin{pmatrix} u_1 & u_2 \end{pmatrix}}_U \underbrace{\begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}}_\Sigma \underbrace{\begin{pmatrix} v_1^T \\ v_2^T \end{pmatrix}}_{V^T} \tag{1.11}$$

### 1.1.3 Application, inference of signals

It is said that SVD can pick out "interesting" signals from data. I'll illustrate this with a simple toy model. You have a rank-1 correction to a matrix full of noise, you want to infer this correction & it's strength. This is called the Spiked Wigner model.

$$A = \lambda x x^T + W \tag{1.12}$$

where  $W \sim \text{GOE}(0, \Sigma)$ , this means that  $\mathbb{E}[W] = 0$  and  $\mathbb{E}[W^T W] = \sigma^2 \mathbb{I}$ .

$$\mathbb{E}[A^T A] = \lambda^2 x x^T x x^T + \sigma^2 \mathbb{I} \quad (1.13)$$

$$= \lambda^2 x^2 x x^T + \sigma^2 \mathbb{I} \quad (1.14)$$

The eigensystem for this is

- One eigenvector  $x$ , with eigenvalue  $\lambda^2 \|x\|^4 + \sigma^2$
- $d - 1$  eigenvectors  $v$  s.t.  $v \perp x$ , with eigenvalue  $\sigma^2$ .

Recall the SVD  $A = U \Sigma V^T$ . the  $V$  were the eigenvectors of  $A^T A$ , which we found contains the signal we wanted to infer. The  $\Sigma$  contains the eigenvalues of  $A^T A$ , which contain information on the strength of the signal  $\lambda$  and noise  $\sigma^2$ .

For more technical readers, this is just a heuristic illustration. Don't forget about the BBP phase transition.

## 1.2 Matrix Calculus

In the next week you'll have to optimize your neural network. Optimization scheme rely on gradient access to your target function, so you'll have to learn how to compute derivatives of vectors & such.

A tiny note on notation. Consider a vector  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^{d \times 1}$  (for example a trajectory), where  $\mathbf{x} = (x_1, \dots, x_d)^T$  :

$$\frac{d\mathbf{x}}{dt} \equiv \begin{pmatrix} \frac{\partial x_1}{\partial t} \\ \vdots \\ \frac{\partial x_d}{\partial t} \end{pmatrix}. \quad (1.15)$$

Now consider the scalar field  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\frac{\partial \phi}{\partial \mathbf{x}} \equiv \begin{pmatrix} \frac{\partial \phi}{\partial x_1} & \dots & \frac{\partial \phi}{\partial x_d} \end{pmatrix} = (\nabla \phi)^T \in \mathbb{R}^{1 \times d}. \quad (1.16)$$

The notation is written s.t.  $\frac{\partial \phi}{\partial \mathbf{x}} \mathbf{x}$  is a sensible matrix operation, and yields a scalar. Now consider the vector field  $\mathbf{y} : \mathbb{R}^d \rightarrow \mathbb{R}^n$  (for example change of coordinates)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \equiv \begin{pmatrix} - & \frac{\partial y_1}{\partial \mathbf{x}} & - \\ & \vdots & \\ - & \frac{\partial y_n}{\partial \mathbf{x}} & - \end{pmatrix} \in \mathbb{R}^{n \times d} \quad (1.17)$$

this is also known as the Jacobian. The notation is written s.t.  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \mathbf{x}$  is a sensible matrix multiplication.

Finally consider a matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$

$$\frac{\partial \mathbf{M}}{\partial t} \equiv \begin{pmatrix} \frac{\partial M_{11}}{\partial t} & \cdots & \frac{\partial M_{1n}}{\partial t} \\ \vdots & \ddots & \vdots \\ \frac{\partial M_{m1}}{\partial t} & \cdots & \frac{\partial M_{mn}}{\partial t} \end{pmatrix} \quad (1.18)$$

$$\frac{\partial t}{\partial \mathbf{M}} = \begin{pmatrix} \frac{\partial t}{\partial M_{11}} & \cdots & \frac{\partial t}{\partial M_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial t}{\partial M_{m1}} & \cdots & \frac{\partial t}{\partial M_{mn}} \end{pmatrix} \quad (1.19)$$

### Problem 1

Derive the back-propagation for a two layer neural network. That is given

$$\mathcal{L} = (y - \hat{y})^2 \quad (1.20)$$

$$\hat{y} = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\mathbf{w}_i^T \mathbf{x} + b_i) \quad (1.21)$$

where  $y, \hat{y}, a_i, b_i \in \mathbb{R}$  are scalars, and  $\mathbf{w}_i, \mathbf{x} \in \mathbb{R}^d$  are vectors. Note  $\mathbf{w}_i$  is not the entry of a vector, there are  $i = 1, \dots, N$   $\mathbf{w}_i$  vectors.

Compute

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_i} \quad (1.22)$$

this is the notation for the gradient w.r.t.  $\mathbf{w}_i$ .

### 1.2.1 Matrix Inversions

**Fact 2 (Sherman-Morrison)** *Let  $A$  be an invertible square matrix, and  $u, v$  be vectors.*

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \quad (1.23)$$

*Proof:* Here's a constructive proof. Say you want to solve for  $x$

$$(A + uv^T)x = y \quad (1.24)$$

$$x = A^{-1}y - A^{-1}uv^T x \quad (1.25)$$

Notice

$$v^T x = v^T A^{-1}y - v^T A^{-1}uv^T x \quad (1.26)$$

$$(1 + v^T A^{-1}u)v^T x = v^T A^{-1}y \quad (1.27)$$

$$v^T x = \frac{v^T A^{-1}}{1 + v^T A^{-1}u} y \quad (1.28)$$

Therefore

$$x = \left( A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) y \quad (1.29)$$

□

The idea is that a rank one perturbation  $uv^T$  to a full rank matrix  $A$ , yields an inverse which is a rank one perturbation to  $A^{-1}$ .

**Fact 3 (Woodbury)** *A generalization of the previous fact is*

$$(A + UCV)^{-1} = A^{-1} + A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1} \quad (1.30)$$

where  $A$  is  $n \times n$ ,  $C$  is  $k \times k$ ,  $U$  is  $n \times k$  and  $V$  is  $k \times n$ .

Proof left as exercise.

### 1.3 Time Complexity

Since we're talking about these things in the context of a computational class, it'll be good to recap the time complexity of such algorithms. Just keep these in the back of your mind.

- Matrix multiplication:  $\mathcal{O}(n^{2.8})$ .
- Matrix inverse implemented in `numpy.linalg.solve`:  $\mathcal{O}(n^3)$ .
- SVD for a  $n \times m$  matrix (s.t.  $n \leq m$ ):  $\mathcal{O}(mn^2)$ .
- Determinant  $\mathcal{O}(n^3)$

As a final note, the time complexity of an algorithm doesn't translate to the actual run time of an algorithm.

See [https://en.wikipedia.org/wiki/Computational\\_complexity\\_of\\_mathematical\\_operations#Matrix\\_algebra](https://en.wikipedia.org/wiki/Computational_complexity_of_mathematical_operations#Matrix_algebra) for more information.

### 1.4 References

## References

- [1] Michael Hutchings, Notes on singular value decomposition for Math 54, <https://math.berkeley.edu/~hutching/teach/54-2017/svd-notes.pdf>.
- [2] Gregory Gundersen, Singular Value Decomposition as Simply as Possible, <https://gregorygundersen.com/blog/2018/12/10/svd/>
- [3] Leslie Lamport (1994) *TEX: a document preparation system*, Addison Wesley, Massachusetts, 2nd ed.

## 2 PyTorch 101

### 2.1 Training Dynamics

Thinking of back-propagation as a high-level chain rule is not enough to get a neural network to train. For example, here's some food for thought...

- We want to  $\arg \min_{\theta} \mathcal{L}(\theta)$ , what scheme do you implement? Note that  $\mathcal{L}(\theta)$  implicitly depends on the data and model architecture, so perhaps that'll affect your answer.
- In the model architecture, how do you initialize the model weights?

### 2.2 Whitening

The following analysis will assume the data distribution  $\mathcal{D} = \{x^{(i)}\}_i \sim^{iid} p_{data}$ , has vanishing  $\mathbb{E}[x^{(i)}] = 0$  and covariance  $\mathbb{E}[x^{(i)}x^{(i)T}] = \sigma^2 \mathbb{I}$ .

You might ask, is it this a fair assumption, and the answer is totally. You can always transform data to have these properties.

Consider the data matrix  $X \in \mathbb{R}^{b \times d}$  (each row is a  $d$ -dimensional data point)

$$X \equiv \begin{pmatrix} - & x^{(1)} & - \\ \vdots & \vdots & \vdots \\ - & x^{(b)} & - \end{pmatrix} \quad (2.1)$$

### 2.3 Back propagation

Using the rules from our matrix cheat sheet. Let  $x \in \mathbb{R}^m$ . Consider a linear layer

$$\text{Linear} : \mathbb{R}^m \rightarrow \mathbb{R}^n \quad (2.2)$$

$$z(x) = \text{Linear}(x) = Wx + b \quad (2.3)$$

$$\frac{\partial}{\partial W_{ij}} \text{Linear}(x) = x_j \quad (2.4)$$

$$\frac{\partial}{\partial W_{ij}} \sigma(z) \quad (2.5)$$

### 2.4 Weight Initialization

As you propagat

### 2.5 PyTorch Tutorial

See `pytorch_tutorial.ipynb` in the repo.



### 3 Review of Probability

**Definition 3 (Conditional Probability)**

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} \quad (3.1)$$

Notice that  $\mathbb{P}[A \cap B] = \mathbb{P}[B \cap A]$ , this allows us to relate  $\mathbb{P}[A|B]$  and  $\mathbb{P}[B|A]$ .

$$\mathbb{P}[A|B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]} \quad (3.2)$$

$$= \frac{\mathbb{P}[B \cap A]}{\mathbb{P}[B]} \quad (3.3)$$

$$\boxed{\mathbb{P}[A|B] = \frac{\mathbb{P}[B|A] \mathbb{P}[A]}{\mathbb{P}[B]}} \quad (3.4)$$

This is **Bayes' Formula**.

**Definition 4 (Probability Density Function)** *A function with the following properties is a **probability density***

- *Positive:*  $p : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$
- *Normalized:*  $\int_{\mathcal{X}} p(x) dx = 1$

*It is interpreted as the probability of observing an event  $A \subset \mathcal{X}$  as*

$$\mathbb{P}[x \in A] = \int_{A \subset \mathcal{X}} p(x) dx \quad (3.5)$$

The nice part of densities is that you can compute statistics with that. I.e. what's the mean, variance.

$$\mathbb{E}_{x \sim p}[f(x)] = \int_{\mathcal{X}} f(x) p(x) dx \quad (3.6)$$

**Definition 5 (Characteristic Function)** *Consider the probability distribution  $p_X$ . It has an associated **characteristic function**  $\varphi_X$  which is its Fourier Transform*

$$\varphi_X(k) = \int_{\mathbb{R}} e^{ikx} p(x) dx = \mathbb{E}_{x \sim p}[e^{ikx}] \quad (3.7)$$

### 4 Review of Statistics & Loss Functions

In machine learning, we adjust a model's parameters  $\theta$  to minimize a loss function  $\mathcal{L}(\theta)$ . There's a bunch so I think it's nice to hear where they come from. We'll cover

- Mean squared error (MSE)

$$\mathcal{L}(\theta) = \sum_{i=1}^n \|y_i - f_{\theta}(x_i)\|^2$$

- Cross entropy

$$\mathcal{L}(\theta) = \sum_{i=1}^n \|\cdot\|$$

- MSE + L2 Regularization (Ridge)

$$\mathcal{L}(\theta) = \sum_{i=1}^n \|y_i - f_{\theta}(x_i)\|_2^2 + \lambda \|\theta\|_2^2$$

## 4.1 Maximum Likelihood Inference & Mean Squared Error

Say you have the dataset  $\mathcal{D} = \{(y_i, x_i)\}_{i=1}^n$  (which we assume you observed in an iid way). You believe that  $y_i$  is a noisy observation of some model  $f_\theta(x_i)$ . Your objective is to come up with the "best" estimate of the parameter  $\theta$  which matches the data  $\mathcal{D}$ ... You think about it for some while, and realize you maximize the probability of seeing the data for a given  $\theta$ . This is **maximum likelihood estimation (MLE)**.

To illustrate this method (and all others), we have to assume a particular model. So let's say you believe the noise is additive & gaussian:

$$y_i = f_\theta(x_i) + \epsilon_i, \text{ where } \epsilon_i \sim_{iid} \mathcal{N}(0, \mathbb{I}) \quad (4.1)$$

Since  $\epsilon_i$  is a random variable, you can interpret  $y_i$  as a random variable as well.

$$y_i \sim \mathcal{N}(f_\theta(x_i), \mathbb{I}) \quad (4.2)$$

$$p(y_i|\theta) \propto \exp\left(-\frac{1}{2}(y_i - f_\theta(x_i))^2\right) \quad (4.3)$$

$$\log p(y_i|\theta) = -\frac{1}{2}(y_i - f_\theta(x_i))^2 + \text{Constant w.r.t. } \theta \quad (4.4)$$

I've wrote the log prob for reasons that will become clear in a moment.

Note you have more data  $\{(y_i, x_i)\}_{i=1}^n$  (which is all iid), so you actually have a joint distribution.

$$p(y_1, \dots, y_n|\theta) = \prod_i p(y_i|\theta) \quad (4.5)$$

We'll call this the **likelihood**  $L(\theta)$  (that is the likeliness / probability of seeing the data given a configuration of model parameters). For MLE, you choose  $\hat{\theta}$  which maximizes the likelihood. However arg max of a product of functions is quite difficult, we can compose the function w/ a monontonic function, and that leaves the arg max invariant.

$$\log L(\theta) = \log p(y_1, \dots, y_n|\theta) \quad (4.6)$$

$$= \sum_i \log p(y_i|\theta) \quad (4.7)$$

$$\propto \sum_i (y_i - f_\theta(x_i))^2 \quad (4.8)$$

This recovers the MSE loss.

## 4.2 Cross Entropy & Another MLE

## 4.3 L2 Regularization

In Bayesian statistics, instead of asking what's the probability of seeing the data given a model parameter, we ask *what's the probability of seeing a model parameter given the data?* We can formalize the inverse question using Bayes' theorem

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)} \quad (4.9)$$

- $p(\theta)$  is your prior. It encodes your prior beliefs into the distribution.

- $p(y|\theta)$  is the likelihood (from the previous sections)
- $p(\theta|y)$  is the posterior. It accounts for your prior beliefs & what the data says (likelihood).
- $p(y)$  is the evidence. I won't say much about it today.

If you ask, what's the parameter maximizes the posterior (probability of seeing a parameter given the data), this is called **maximum a posteriori estimation (MAP)**.

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta|y) \quad (4.10)$$

For an example, let's assume we have the additive noise model

$$y_i = f_{\theta}(x_i) + \epsilon_i \quad (4.11)$$

and that you believe the weights should look distributed according to a Gaussian

$$p(\theta) = \mathcal{N}(0, \lambda^{-1}\mathbb{I}) \quad (4.12)$$

You can see that log posterior has the form

$$\log p(\theta|y) = \sum_i \|y_i - f_{\theta}(x_i)\|^2 + \lambda \|\theta\|^2 \quad (4.13)$$

## 4.4 Minimizing the loss function

- The value of the MSE, in a traditional statistics setting, tells you about the uncertainty quantification of the model. However ML models tend to not obey this.
- Difficulty of optimizing via oracle access.
- However! Do you even want to perfectly minimize the loss function? Memorization.

## 5 Linear Regression

Consider making iid noisy observations of data  $\{(x_i, y_i)\}_{i=1}^n$ , where  $x_i \in \mathbb{R}^p$  and  $y_i \in \mathbb{R}$ . We'll assume that the noise is additive, that is

$$y_i = f(x_i) + \epsilon_i \quad (5.1)$$

where  $f(x) = \beta^T x$  is the model (we've assumed it's linear for this discussion) and the noise is gaussian  $\epsilon_i = \mathcal{N}(0, \sigma_i^2)$  (which is another assumption for this discussion). Since  $\epsilon_i$  is a random variable, this implies that  $y_i$  is also a random variable

$$y_i|\beta = \mathcal{N}(y_i; \beta^T x_i, \sigma_i^2) \quad (5.2)$$

This is just one observation, but in fact, we have a joint distribution  $p(y|x) \equiv p(y_1, \dots, y_N|x_1, \dots, x_N)$  over all observations, which we'll call the **likelihood**. Since observations are iid, it factorizes.

$$p(y|\beta) = \prod_i p(y_i|x_i) \quad (5.3)$$

Your task is to find the  $\beta$  which "best" describes the data. I'll note that "best" is subjective and we'll discuss consequences of this later.

## 5.1 Frequentist, Maximum Likelihood Estimator

One method is **maximum likelihood estimation**, that is you select the parameters which is the global maximizer of the likelihood. Why? Just read off what you're doing: adjust  $\beta$  s.t. the probability of having this combination of  $y$ 's (given  $x$ 's) is highest.

Apart from being very intuitive, there are also strong theoretical guarantees (which I won't have time to prove) (Notation: when I generically talk about model parameters, we use  $\theta$ )

- Consistency:  $\lim_{n \rightarrow \infty} \hat{\theta}_n = \theta$
- Normality:  $\hat{\theta}_n \sim \mathcal{N}(0, \mathcal{I})$  (where  $\mathcal{I}$  is the fisher information matrix)
- Efficiency:  $\text{Var}(\hat{\theta}) \geq 1/\mathcal{I}(\theta)$ .

Since  $\arg \max$  is invariant under compositions of monotonic functions, we can maximize the log-likelihood which emits a nicer function

$$\log p(y|x) = \sum_i \log p(y_i|x_i) \quad (5.4)$$

$$= \sum_i \log \mathcal{N}(y_i; \beta^T x_i, \sigma_i^2) \quad (5.5)$$

$$= \sum_i -\frac{1}{2} \frac{(y_i - \beta^T x_i)^2}{\sigma_i^2} + \text{Constant} \quad (5.6)$$

A small comment, this is why you "minimize the squared error" when fitting straight lines in lab, you have been secretly doing maximum likelihood inference this whole time. Notice this is a quadratic form, so you can rewrite it using matrix multiplication

$$\sum_{i=1}^n (y_i - \beta^T x_i)^2 / \sigma_i^2 = (y - X\beta)^T \Sigma^{-1} (y - X\beta) \quad (5.7)$$

$$\text{where: } y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n, \quad (5.8)$$

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^p \quad (5.9)$$

$$\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2) \in \mathbb{R}^{n \times n} \quad (5.10)$$

$$X = \begin{pmatrix} - & x_1^T & - \\ - & x_2^T & - \\ & \vdots & \\ - & x_n^T & - \end{pmatrix} \in \mathbb{R}^{n \times p} \quad (5.11)$$

From here we can find the argmax of the quantity

$$0 = \left. \frac{\partial \log p(y|x)}{\partial \beta} \right|_{\beta=\hat{\beta}} = X^T \Sigma^{-1} (y - X\beta) \quad (5.12)$$

$$\implies X^T \Sigma^{-1} y = X^T \Sigma^{-1} X \hat{\beta} \quad (5.13)$$

$$\boxed{\hat{\beta}_{MLE} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y} \quad (5.14)$$

and find the maximum likelihood estimate for  $\beta$ .

Now we can talk about inference. Say your boss gives you new data  $X_*$ , and you're asked what is the corresponding  $\hat{y}_*$ . You'll report back

$$\boxed{\hat{y}_* = X_* \hat{\beta}_{MLE}} \quad (5.15)$$

## 5.2 Bayesian Linear Regression

In the Bayesian framework, you're asked what is the probability of seeing the model parameters *given* the data  $p(\beta|y)$ . You can calculate this using Bayes's formula

$$p(\beta|y) = \frac{p(y|\beta)p(\beta)}{p(y)} \quad (5.16)$$

# 6 Double Descent

## 6.1 Soft Inductive Biases

Another way to conceptualize this is **soft inductive biases** (see Andrew Gordon Willson's paper <https://arxiv.org/pdf/2503.02113>).

# 7 Training Large Models

## 7.1 Transformers

## 7.2 $\mu$ P Optimizer

# 8 Geometric Deep Learning

## Part I

# Generative Models

## 9 "Old" Generative Models

### 9.1 Variational Autoencoders

### 9.2 Generative Adversarial Networks (GANs)

### 9.3 Denoising Diffusion Probabilistic Models (DDPMs)

## 10 Modern Generative Models

### 10.1 Review of Non-Equilibrium Statistical Mechanics

Consider a classical particle at position  $X_t$  at time  $t$  (usually we use the notation  $x(t)$ , but I'll use  $X_t$  as it's the modern stochastic calculus notation) moving in a potential  $V(x)$ . The

weird thing is this particle appears jittery, it feels a bunch of random forces due to thermal fluctuations

$$m\ddot{X}_t = -V(X_t) + \xi_t \quad (10.1)$$

## 10.2 Measure Transport

**Theorem 2 (Fokker-Planck Equation)** *Consider a stochastic process*

$$dX_t = \mu_t(X_t) dt + \sigma_t(X_t) dW_t \quad (10.2)$$

$$X_0 \sim p_{base} \quad (10.3)$$

*The stochastic process emits a probability distribution at every points in time (notationally  $X_t \sim p_t$ ), where the distribution  $p_t$  satisfies a partial differential equation called the **Fokker-Planck equation***

$$\partial_t p_t(x) = -\nabla \cdot (\mu_t(x) p_t(x)) + \frac{1}{2} \sigma_t^2(x) \Delta(p_t(x)) \quad (10.4)$$

$$p_{t=0}(x) = p_{base} \quad (\text{Boundary condition}) \quad (10.5)$$

A small remark, take  $\sigma \rightarrow 0$  and you recover the transport equation

**Theorem 3 (Transport Equation)** *Consider the deterministic process*

$$dX_t = \mu(X_t) dt \quad (10.6)$$

$$X_0 \sim p_{base} \quad (10.7)$$

*The associated probability distribution  $X_t \sim p_t$  satisfies the partial differential equation called the **transport equation***

$$\partial_t p_t(x) = -\nabla \cdot (\mu_t(X_t) p_t(x)) \quad (10.8)$$

This increases the design space.

## 10.3 Score Based Diffusion

## 10.4 Flow Matching

## 10.5 Stochastic Interpolants