

mlfp (Python 3.12.0)

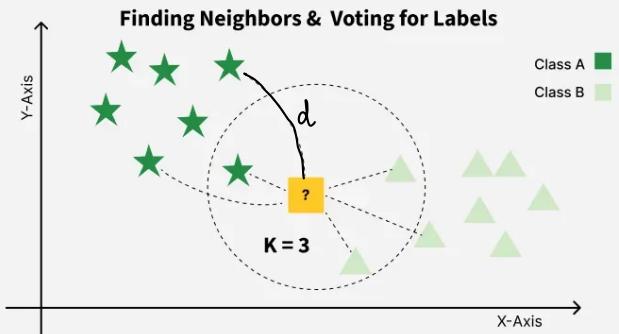
Model is not actually doing general learning, but simply = memorizing "the dataset \Rightarrow Overfitting

\Rightarrow Dataset too small (not enough for general learning) & Model has more learning parameters than it needs to correctly model the dataset.

1

k-NN Algorithm:

i) Euclidean Distance $d(x, x_i) = \|x - x_i\| = \sqrt{\sum_{j=1}^{d-\text{dim.}} (x_j - x_{ij})^2}$ Minkowski: $d = \|x - x'\|_p$



\Rightarrow no assumptions about the underlying data distribution it makes a non-parametric and distance-based learning method.

→ Data has lots of noise or outliers, using a large k making prediction more stable.

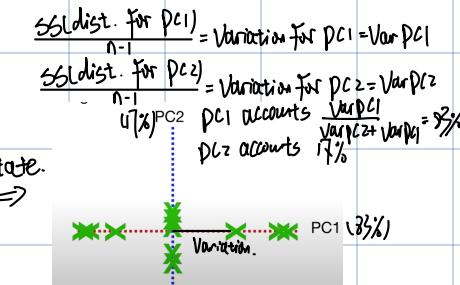
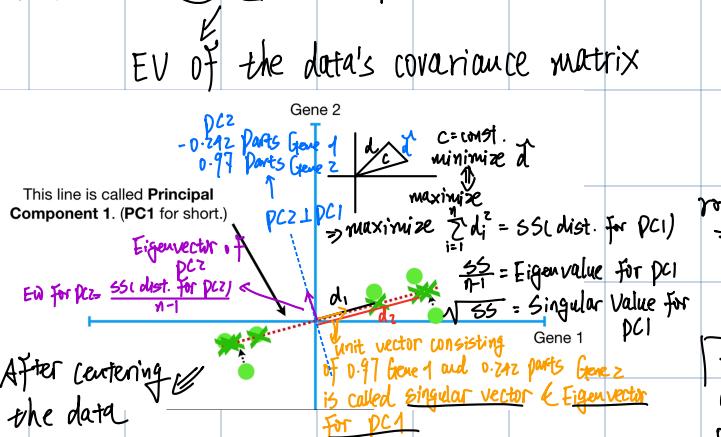
\Rightarrow k too large \Rightarrow model become simple, loss of important patterns
 $\quad \quad \quad$ = underfitting"

$$\text{Manhattan: } d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

Selecting k : Cross-Validation, Elbow Method, $\text{RE} \approx N_0 + d$.

② Dimensionality reduction: reduce high-dim Data, lowering computational costs with minimal loss of Info.

I) PCA (Principal component analysis) (linear method)

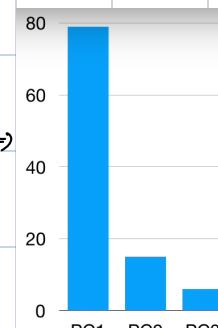
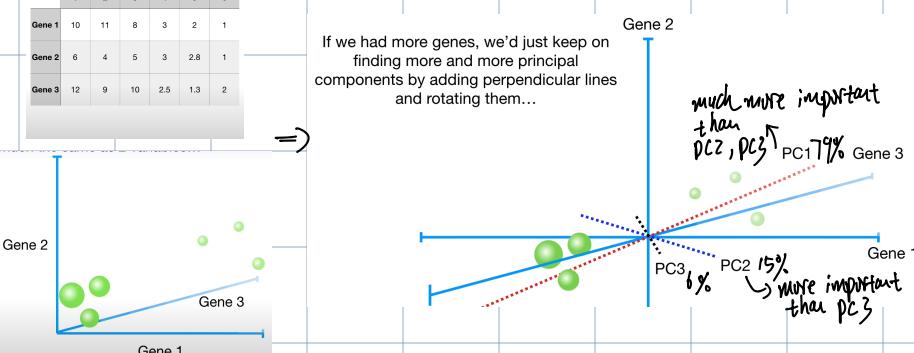


The PC of a collection of points are a sequence $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_P$ of P unit vectors, where the i -th vector is the direction of a line that best fits the data while being \perp to the first $i-1$ vectors. The best-fitting line is defined as one that minimizes the average squared perpendicular distance from the points to the line.

More Variables: i) center data ii) best fitting through origin: PC1 iii) draw PC2, PC3 \perp to each other iv) compute Variation in each PC axis

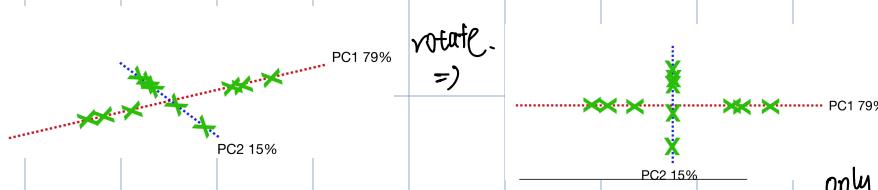
	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1
Gene 3	12	9	10	2.5	1.3	2

If we had more genes, we'd just keep on finding more and more principal components by adding perpendicular lines and rotating them...



only with PC1 PC2
account for 94% of the
variation of data.
(Approx. with 2D, PC1 and
PC2)

v) proj. into
PC1 - PC2 plane.
=>



High-Dim. Data
reduced into 2-dim.
PC1 x PC2 Diagram.

only works well if the first 2 PCs
account for most of the variation in data

capture local structure

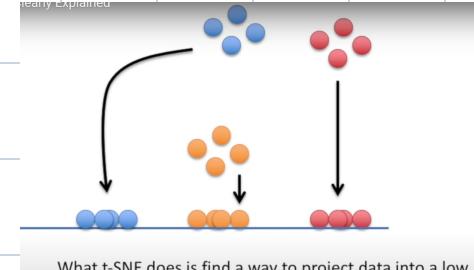
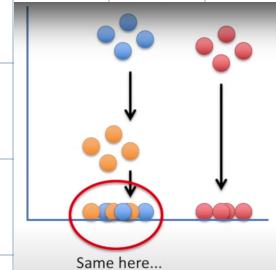
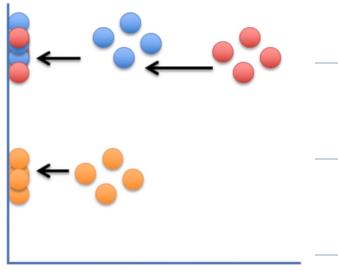
II t-SNE (t-distributed stochastic neighbor embedding) (better als PCA)

PCA works better on linear data (not so good if scatter plot uniform distributed)

=> map high-dim data points into a lower-dim. space, in a way that preserves the local relationships between points

(A) measure similarity between data points in the high-dim. space and representing this similarity as probabilities

(B) Then constructs a similar probability distribution in lower-dim. space and minimizes the difference between the two distributions using a technique called gradient descent



NOTE: If we just projected the data onto one of the axes, we'd just get a big mess that doesn't preserve the original clustering.

What t-SNE does is find a way to project data into a low dimensional space (in this case, the 1-D number line) so that the clustering in the high dimensional space (in this case, the 2-D scatter plot) is preserved.

process

Quest: t-SNE, Clearly Explained

First, measure the distance between two points...

Then plot that distance on a normal curve that is centered on the point of interest...

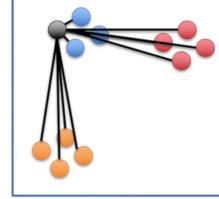
using normal distribu:
distinct points have
very low similarity values
and close points high
similarity.

...lastly, draw a line from the point to the curve. The length of that line is the "unscaled similarity".

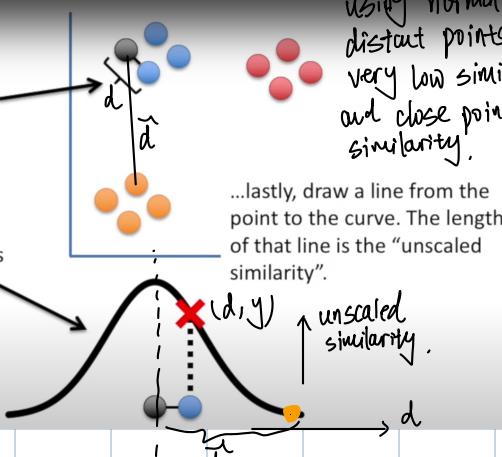
=>

Ultimately, we measure the distances between all of the points and the point of interest...

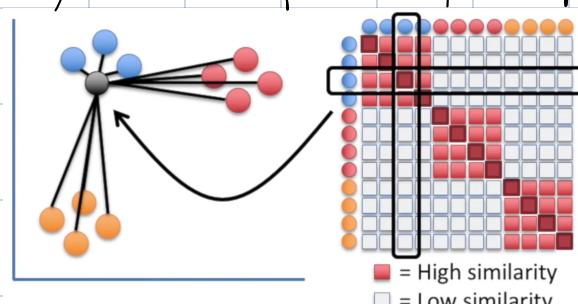
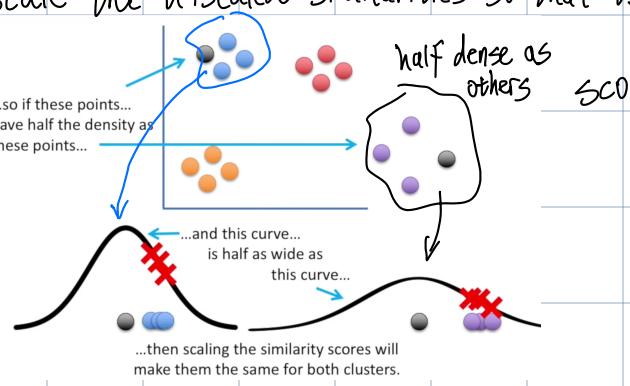
Plot them on the normal curve...



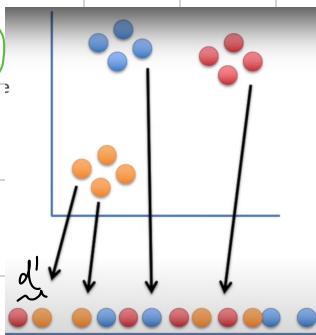
...and then measure the distances from the points to the curve to get the unscaled similarity scores with respect to the point of interest.



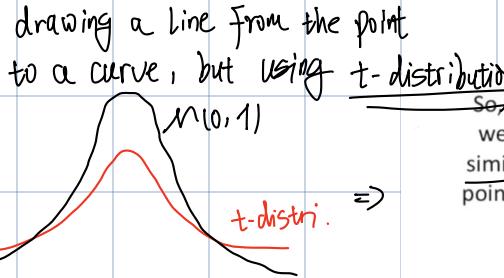
=> Scale the unscaled similarities so that they add to 1 (density cluster, same scaled similarity scores) => do this procedure for all points of interest:



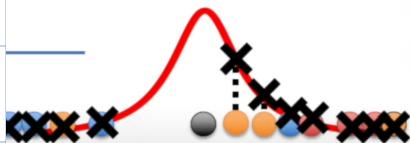
■ = High similarity
□ = Low similarity



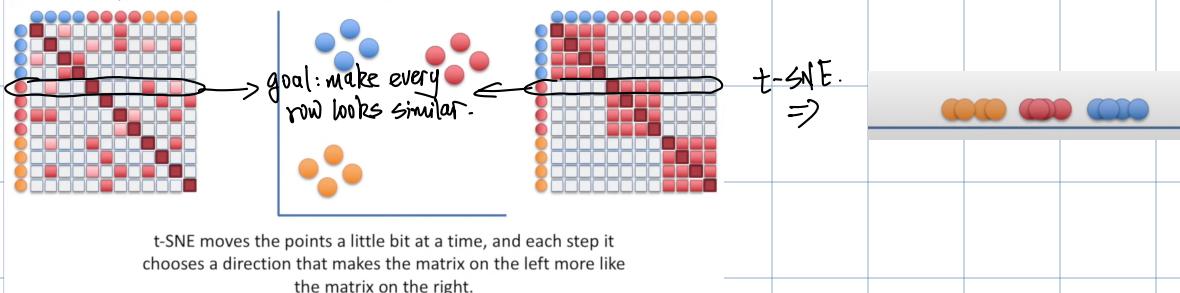
randomly project the data onto the number line, and calculate similarity scores for the points on the number line



So, using a t-distribution, we calculate "unscaled" similarity scores for all the points and then scale them like before.



\Rightarrow scaling \Rightarrow do for all points, end up with matrix of similarity scores on number line.

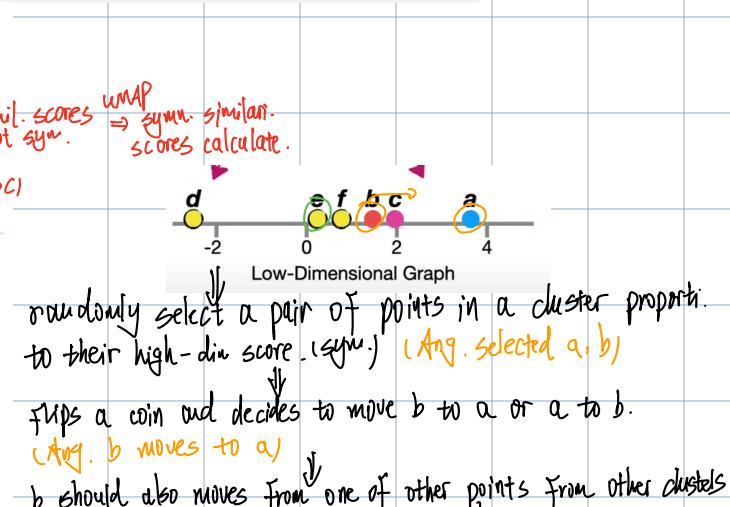
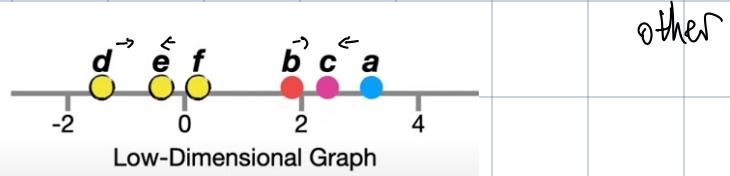
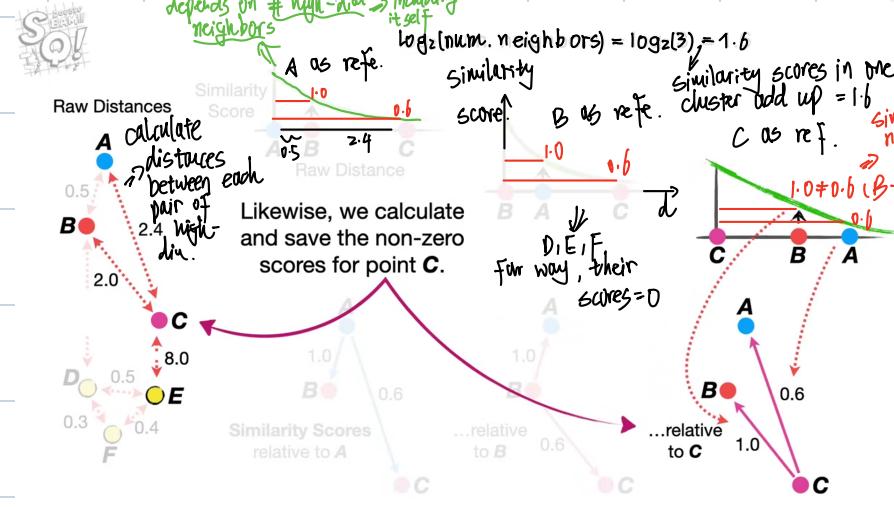
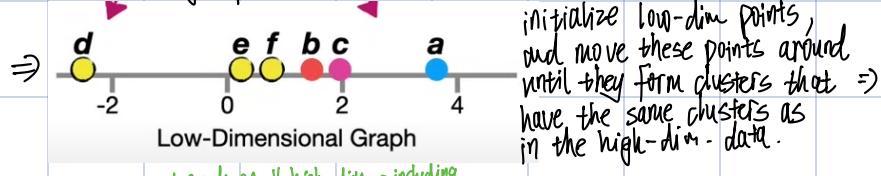


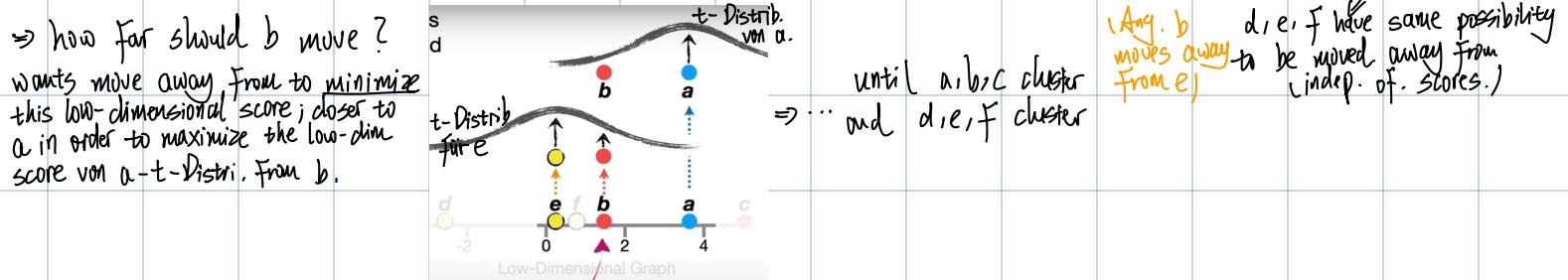
uniform Manifold Approx. and Projection

III UMAP t-SNE's performance suffers with large datasets and using it correctly can be increased speed. + better preservation of data's global structure challenging.

but both algorithms performs strong local structure and group similar categories together; UMAP separates much more clearly the groups of similar categories from each other.

Process: A, B, C seem to form a cluster.; same as D, E, F; Goal of UMAP is to create a low-dim. graph of this data that preserves these high-dim clusters and their relationship to each other





Differences t-SNE & UMAP: ① t-SNE starts with random initialization of low-dim graph ; UMAP uses Spectral Embedding and starts always with the exact same low-dim graph. ② t-SNE moves every single point a little bit by each iteration ; UMAP can move just one point, or a small subset of points , each time and this helps it scale well with super big datasets

③ Advanced Machine Learning.

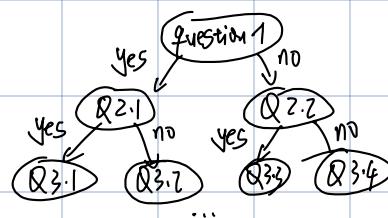
- **Decision Trees** : make decisions (yes or no) by asking a series of questions
 - simple to understand, versatility (can handle both categorical and numerical data), no need for data preprocessing, not affected by outliers or missing values ; Feature selection (using first the most informative features , can be particularly beneficial when dealing large number datasets).
 - Overfitting (can create a tree perfectly classif. training data but performs poorly on test data) ; Bias towards features with more levels. ; Instability (sensitive to small changes in the data, slight variation leads to drastically different trees) ; Difficulty in capturing complex relationships (only yes, no).

• Random Forest.

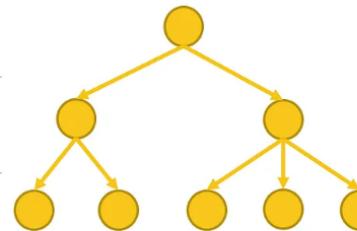
In a Random Forest, instead of just one decision tree making all the decisions, we create an entire “forest” of decision trees. Each tree gives its “opinion” or prediction based on the data it has seen. The final output is then determined by considering the output of all the trees in the forest.

- “Opinion gathering” more exciting by giving each tree a slightly differ. set of data points and features to learn from. (bagging and feature randomness)
- less possible overfitting than single decision Trees. ; handles large Data Set and Feature Spaces

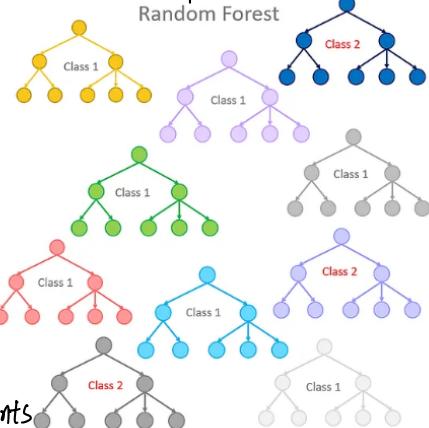
Training of individual trees can be done in parallel ; can provide insights into which features are most import.



Single Decision Tree



team of decision trees

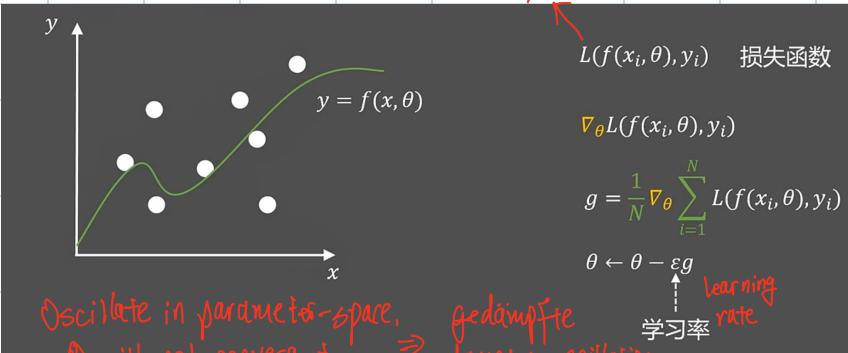


in making predictions.

→ Complexity ; Less Interpretable than single decisi. Tree ; longer prediction time

• Gradient Descent.

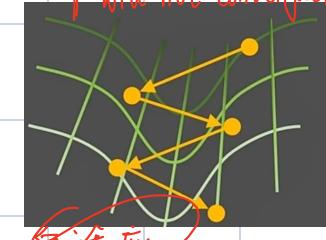
Loss Function minimieren



⇒ random gradient descent : every $m \leq N$ different
 $g = \frac{1}{m} \nabla_\theta \sum_{i=1}^m L(f(x_i, \theta), y_i)$ data points
 (Stochastic gradient descent) (if std. derivation small, although noise but proper convergent Dachsturm).

Oscillate in parameter-space.
↑ will not converge!

gradient descent
damping oscillation.



= Stochastic Gradient Descent

AdaGrad (2011)

$$g = \frac{1}{m} \nabla_\theta \sum_{i=1}^m L(f(x_i, \theta), y_i)$$

$r \leftarrow r + g^2$ value

quickly find convergent direction

$$\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{r} + \delta} g$$

fluctuates, learning rate decreases fast (at beginning); after a. g small

δ is a small number, stable numerical calculation

careful adjust learning rate to avoid loss function oscillate.

RMSProp (2012)

$$g = \frac{1}{m} \nabla_\theta \sum_{i=1}^m L(f(x_i, \theta), y_i)$$

wenn gradient g

$$\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{r} + \delta} g$$

rate decreases fast (at beginning); after a. g small

δ is a small number, stable numerical calculation

to avoid learning rate decreases too fast.

(manually control)

Adam (2014) (impuls + self-adjusting)

$$g = \frac{1}{m} \nabla_\theta \sum_{i=1}^m L(f(x_i, \theta), y_i)$$

control impuls with ρ_1

$$s \leftarrow \rho_1 s + (1 - \rho_1) g$$

control self-adaptation with ρ_2 .

$$r \leftarrow \rho_2 r + (1 - \rho_2) g^2$$

$$\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$$

$$\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$$

$$\theta \leftarrow \theta - \frac{\epsilon \hat{s}}{\sqrt{\hat{r}} + \delta} g$$

• Gradient Boosting Trees.

Learning Rate

Boosted Ensemble = First Tree + η^* Second Tree Loss(Boosted Ensemble) < Loss(First Tree)

→ $F(z) = F(1) + \eta^* \text{second tree}$, second tree = $- \frac{\partial L}{\partial F(1)} = - \frac{\partial \text{Loss Function}}{\partial \text{Previous Model's Output}}$.

→ $F(m) = F(m-1) - \eta^* \frac{\partial L}{\partial F(m-1)}$ (Overfitting) ⇒ Decision Trees are not independent anymore!

• Neural Networks.

① MLP (Multi-Layer-Perceptron)

composition of a linear mapping followed by a non-linear function

$$f(x) = \sigma(\tilde{w}x + b) = \sigma\left(\sum_{i=1}^{n_{out}} (w_i x_i + b_i)\right)$$

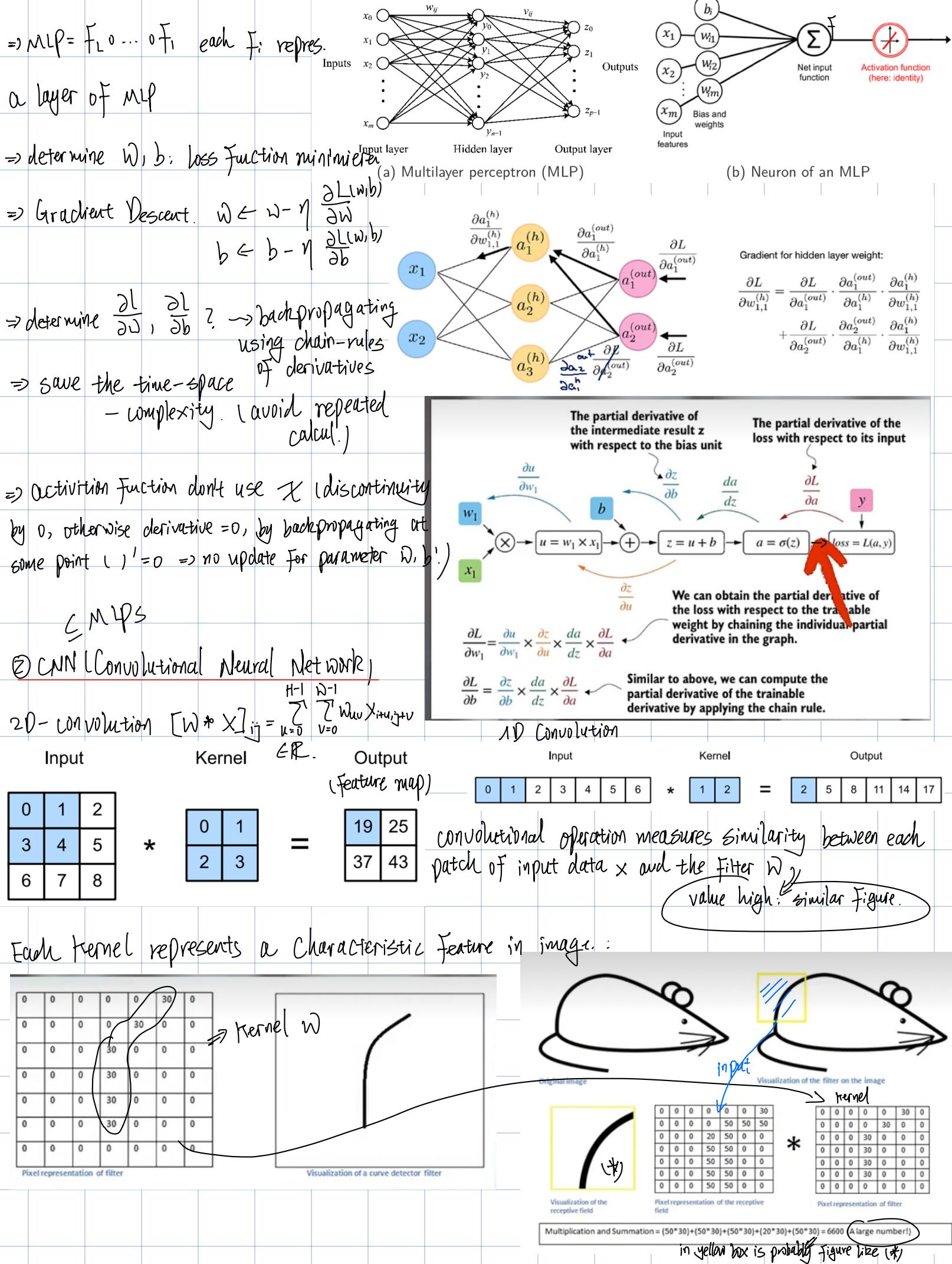
weight. bias activation function

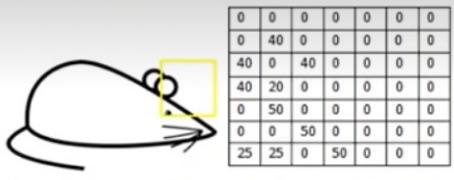
$$x \in \mathbb{R}^{n_{in}} = \mathbb{R}^{n_{in}}$$

$$w \in \mathbb{R}^{n_{out} \times n_{in}}$$

$$b \in \mathbb{R}^{n_{out}}$$

⇒ single Block of MLP





Visualization of the filter on the image

0	0	0	0	0	0	0
0	40	0	0	0	0	0
40	0	40	0	0	0	0
40	20	0	0	0	0	0
0	50	0	0	0	0	0
0	0	50	0	0	0	0
25	25	0	50	0	0	0

Pixel representation of receptive field

*

0	0	0	0	0	30	0
0	0	0	0	30	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	30	0	0	0
0	0	0	0	0	0	0

Pixel representation of filter

Multiplication and Summation = 0 () *this local structure is not same as ()*

