

DIMENSIONALITY REDUCTION AND VISUALIZATION

The curse of dimensionality



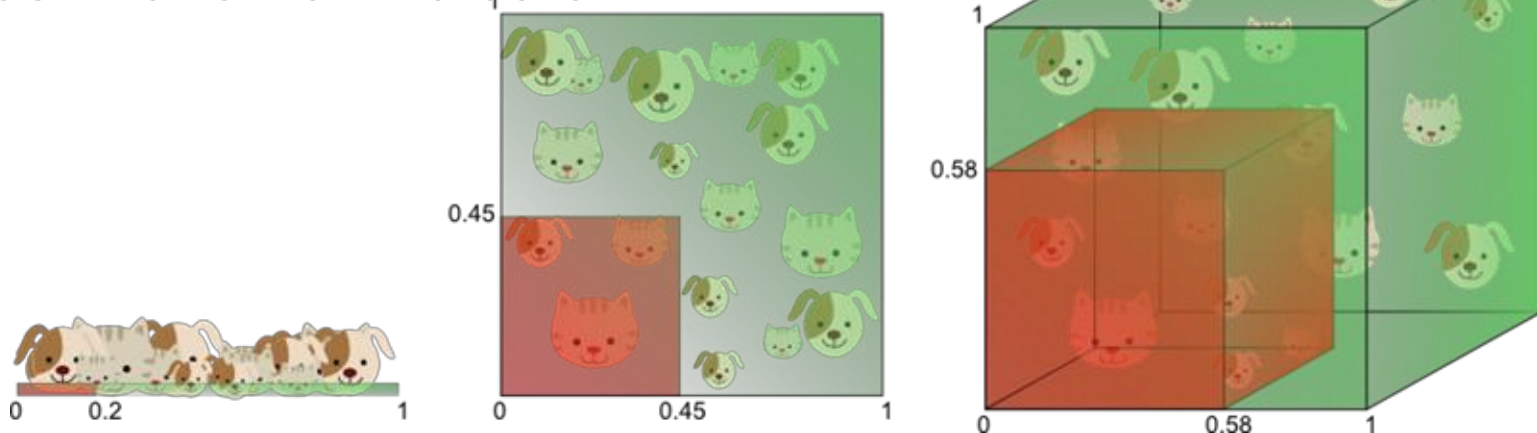
The Curse of Dimensionality

- Harder to visualize or see structure of
 - Verifying that data come from a straight line/plane needs $n+1$ data points
- Hard to search in high dimension – More runtime
- Need more data to get a good estimation of the data

Wind in the morning $X \in \{\text{Calm, Windy}\}$
 PM2.5 level in the afternoon $Y \in \{\text{Low, Med, High}\}$
 PM2.5 level in the evening $Z \in \{\text{Low, Med, High}\}$
 $\text{argmax } P(Z | Y, X) = \text{argmax } P(Y, X | Z) P(Z)$

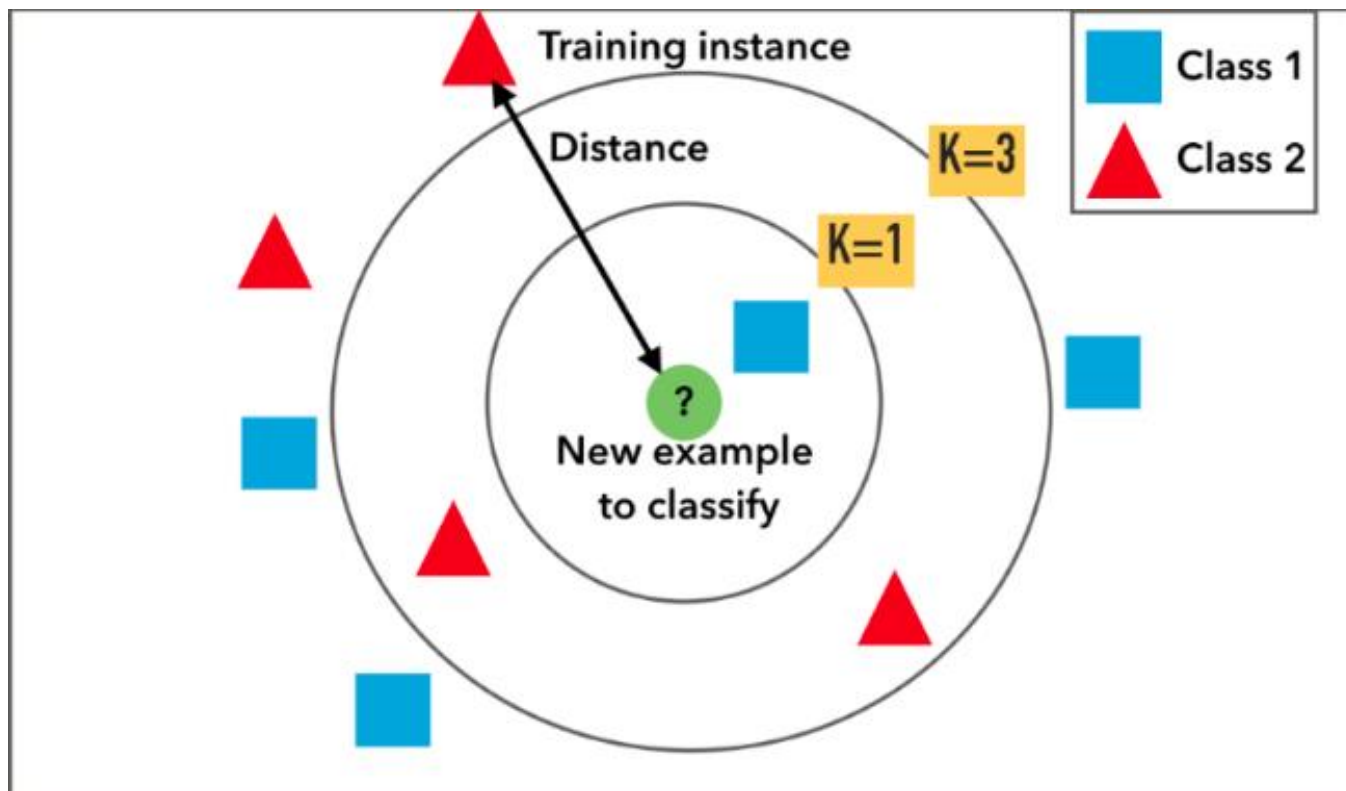
Day	X	Y	Z
1	W	L	M
2	C	M	M
3	W	H	M
4	W	M	H
5	C	M	L
6	W	M	L
7	C	L	H
8	W	H	L

count(Z,Y,X)	Z=L	Z=M	Z=H
X=W,Y=L	0	1	0
X=W,Y=M	1	0	1
X=W,Y=H	1	1	0
X=C,Y=L	0	0	1
X=C,Y=M	1	1	0
X=C,Y=H	0	0	0

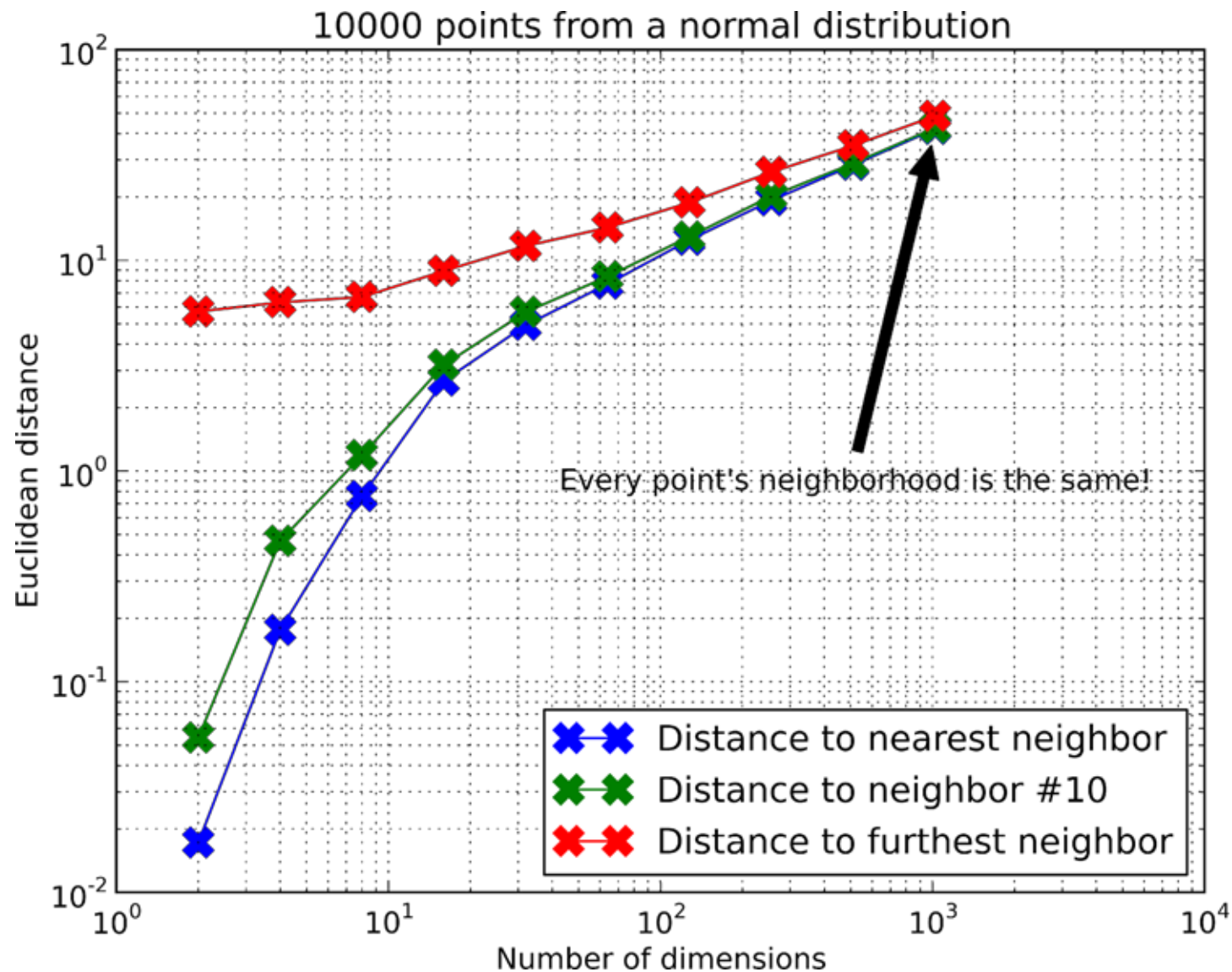


K-Nearest Neighbor Classifier

- Nearest neighbor is susceptible to label noise
- Use the k-nearest neighbors as the classification decision
 - Use majority vote



What's wrong with knn in high dimension?



Combating the curse of dimensionality

- Feature selection
 - Keep only “Good” features
- Feature transformation (Feature extraction)
 - Transform the original features into a smaller set of features

Feature selection vs Feature transform

- Keep original features
 - Useful for when the user wants to know which feature matters
 - Hard to select good features automatically
- New features (a combination of old features)
- Usually more powerful
 - Harder to interpret the model

Feature selection

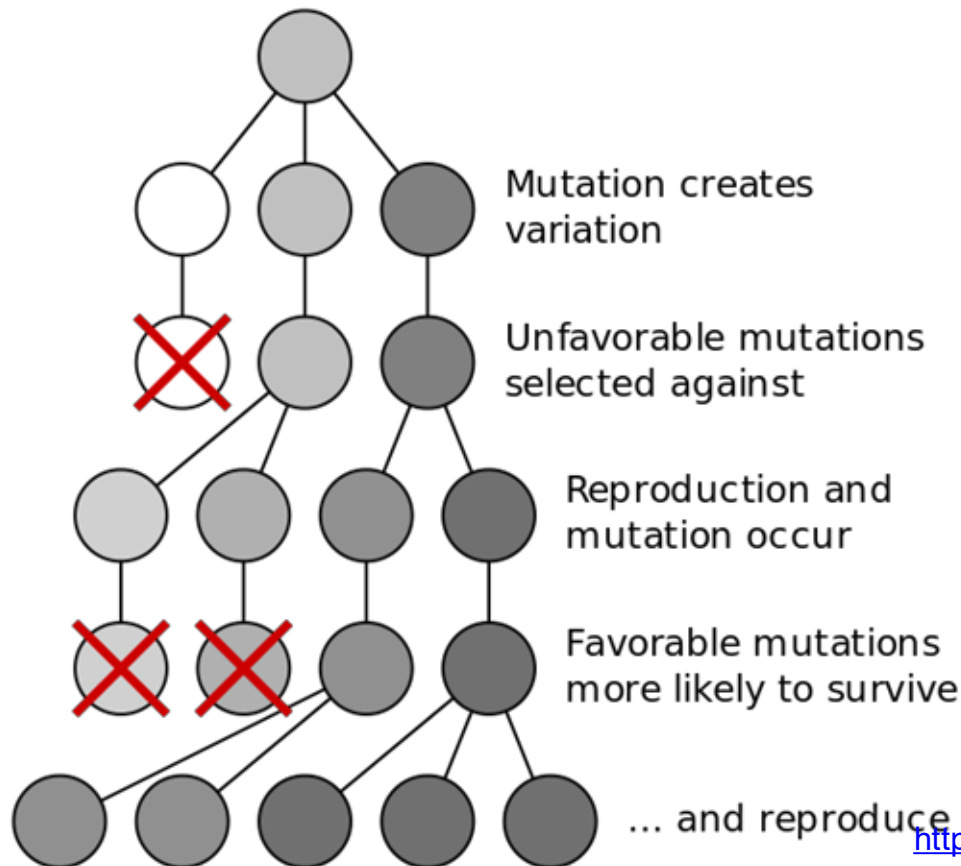
- Hackathon level (time limit days-a week)
 - Drop missing features
 - Low variance column
 - A feature that is a constant is useless. Tricky in practice
 - Forward or backward feature elimination
 - Greedy algorithm: create a simple classifier with $n-1$ features, n times. Find which one has the best accuracy, drop that feature. Repeat.

Feature selection

- Proper methods
 - Algorithm that handles high dimension well and do selection as a by product
 - Regression with L1 regularization
 - Tree-based classifiers: random forest, XGBoost
 - Genetic Algorithm

Genetic Algorithm

- A method based inspired by natural selection
 - No theoretical guarantees but often work

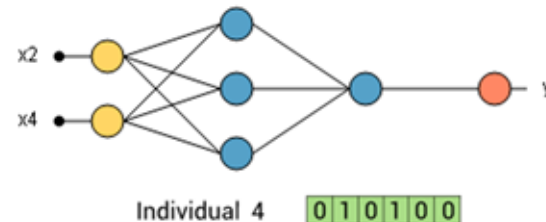
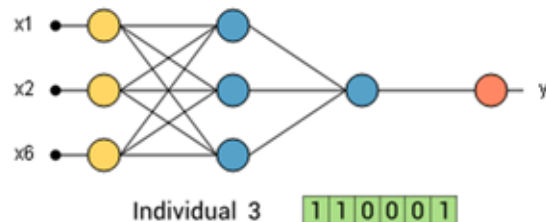
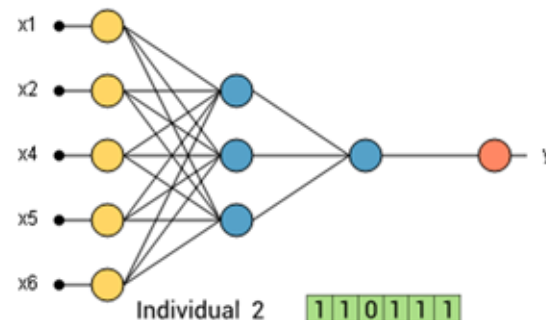
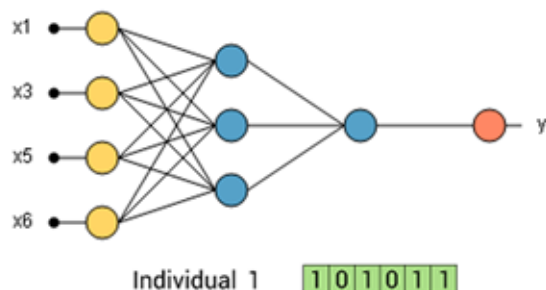


Genetic Algorithm

- Initialization
 - Create N classifiers, each using different subset of features
- Selection process
 - Rank the N classifiers according to some criterion, kill the lower half
- Crossover
 - The remaining classifier breeds offsprings by selecting traits from the parents
- Mutation
 - The offsprings can have mutations by random in order to generate diversity
- Repeat till satisfied

Initialization

- Create N classifiers
- Randomly select a subset of features to use



Selection process

- Score the classifiers and kill the lower half (the amount to kill is also a parameter)

	Selection error	Rank
Individual 1	0.9	1
Individual 2	0.6	3
Individual 3	0.7	2
Individual 4	0.5	4

Crossover

- Breed offsprings by randomly select genes from parents

Individual 3

1	1	0	0	0	1
---	---	---	---	---	---

Individual 4

0	1	0	1	0	0
---	---	---	---	---	---

Offspring 1

0	1	0	1	0	1
---	---	---	---	---	---

Offspring 2

1	1	0	1	0	1
---	---	---	---	---	---

Offspring 3

0	1	0	1	0	1
---	---	---	---	---	---

Offspring 4

1	1	0	0	0	0
---	---	---	---	---	---

Mutation

- Offspring can mutate with some probability to introduce diversity
- Mutation rate is usually $1/k$ where k is the number of features.
 - On average you mutate once per individual

Offspring1: Original

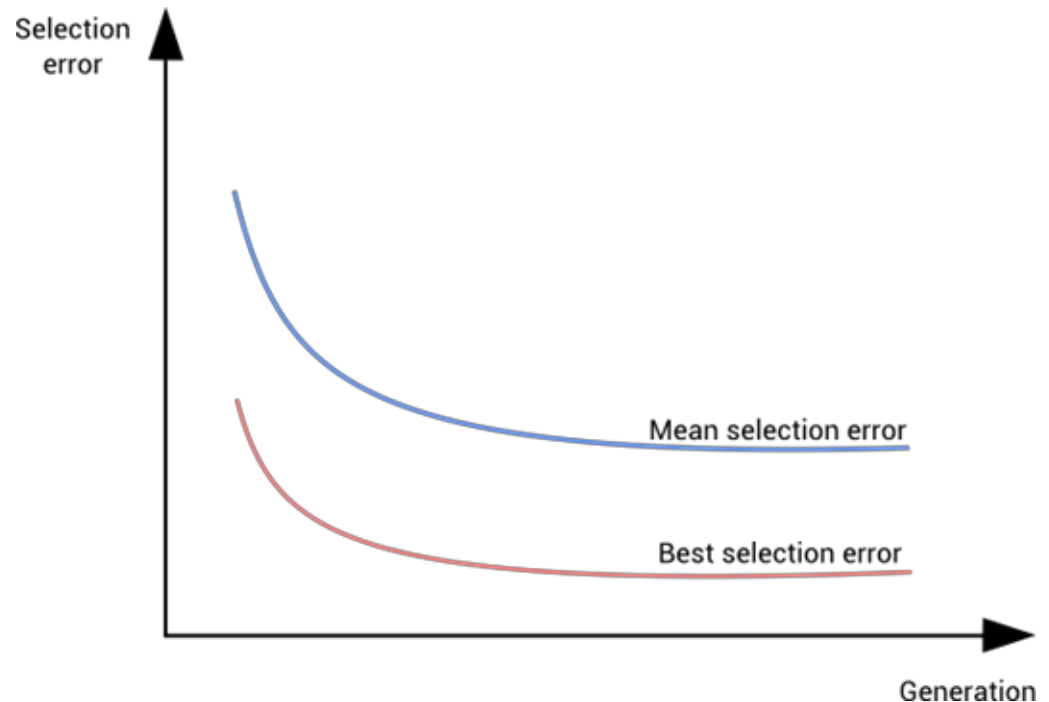
0	1	0	1	0	1
---	---	---	---	---	---

Offspring1: Mutated

0	1	0	0	0	0
---	---	---	---	---	---

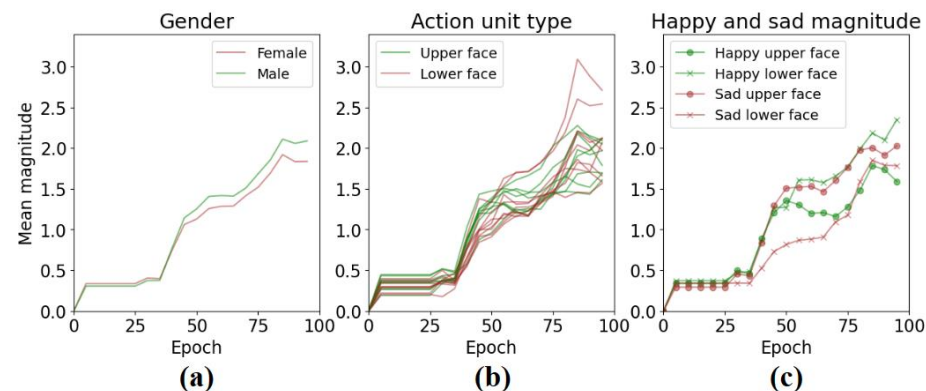
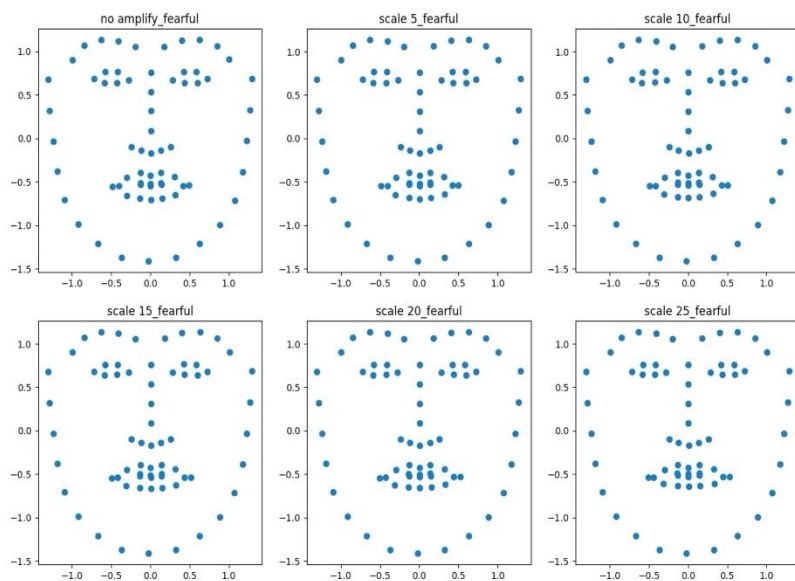
Performance

- Usually performs well. The general population usually gets better (mean). The best performing (individual) also gets better after each generation



Other examples of genetic algorithm

- Tuning hyperparameters in a neural network
 - <https://blog.coast.ai/lets-evolve-a-neural-network-with-a-genetic-algorithm-code-included-8809bece164>
- Tune augmentation algorithms



Generation ->

Feature transformation

- Principal Component Analysis
- Linear Discriminant Analysis (NOT Latent Dirichlet Allocation)
- Random Projections

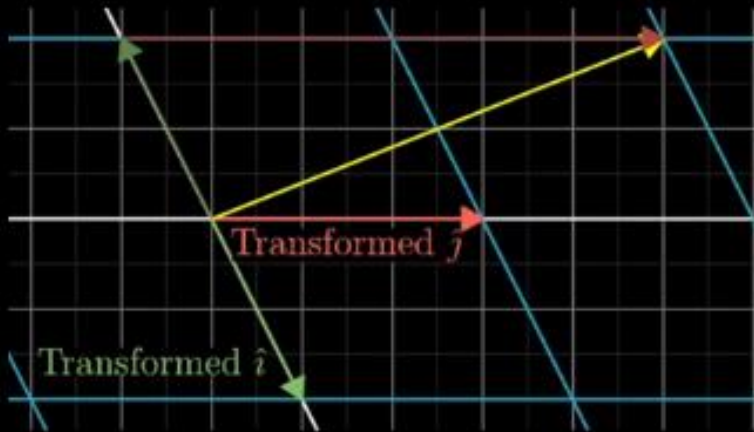
Linear Algebra Review

- Matrix as a sequence of column vectors

“2x2 Matrix”

$$\begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix}$$

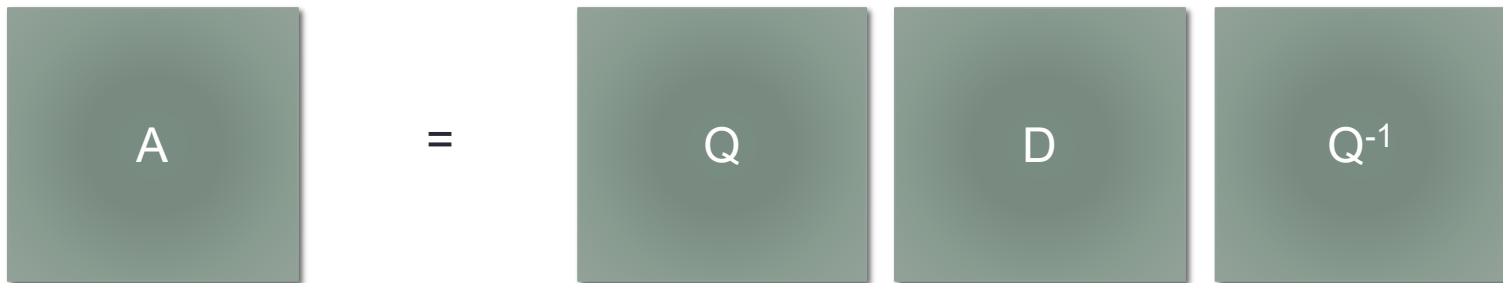
$$\begin{bmatrix} 5 \\ 7 \end{bmatrix}$$



$$5 \begin{bmatrix} 3 \\ -2 \end{bmatrix} + 7 \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

Linear Algebra Review

- View Eigendecomposition (ED) and Singular Value Decomposition (SVD) as rotations and stretches

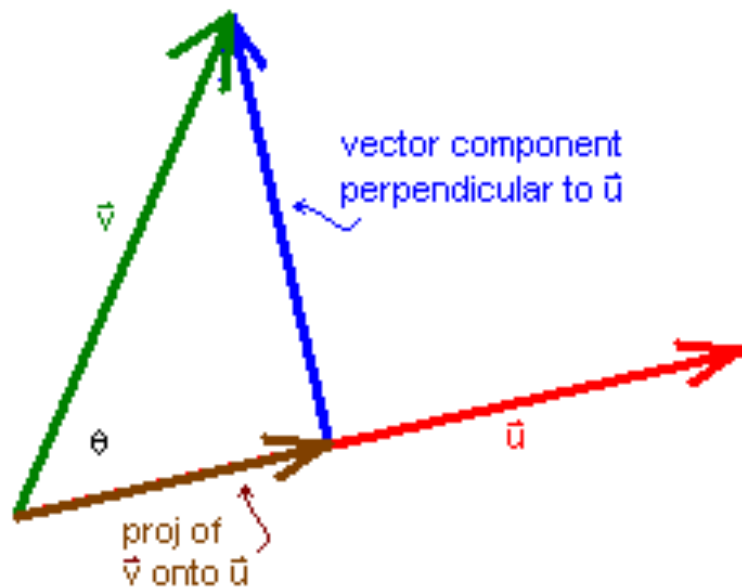

$$A = Q D Q^{-1}$$


$$q_1 \quad q_2 \quad q_3$$

D has eigenvalues on the diagonal
Q is a matrix where i^{th} column is the q^{th} eigenvector

Linear Algebra Review

- Projection as a change of basis
- Change basis from x, y coordinates to be on u



$$u^t v$$

Covariance matrix

- Given a set of RVs, $X_1 X_2 \dots X_n$
- The covariance matrix is a matrix which has the covariance of the i and j RV in position (i,j)

$$\Sigma = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}.$$

Covariance matrix (vector view)

- Given a set of datapoints $x_1 \ x_2 \ \dots \ x_n$
- The covariance matrix based on the data can be written as

$$E[\underset{\substack{\text{vector} \\ \text{of size } n}}{x} \underset{\substack{\text{vector} \\ \text{of size } n}}{x}^T] = [\underset{\substack{\text{matrix} \\ \text{of size } n \times n}}{\square} + \underset{\substack{\text{matrix} \\ \text{of size } n \times n}}{\square} + \dots + \underset{\substack{\text{matrix} \\ \text{of size } n \times n}}{\square}] / n$$

$$\Sigma = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}.$$

Positive semi-definite and covariance matrixes

Definitions for real matrices [\[edit\]](#)

An $n \times n$ symmetric real matrix M is said to be **positive-definite** if $\mathbf{x}^T M \mathbf{x} > 0$ for all non-zero \mathbf{x} in \mathbb{R}^n . Formally,

$$M \text{ positive-definite} \iff \mathbf{x}^T M \mathbf{x} > 0 \text{ for all } \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$$

An $n \times n$ symmetric real matrix M is said to be **positive-semidefinite** or **non-negative-definite** if $\mathbf{x}^T M \mathbf{x} \geq 0$ for all \mathbf{x} in \mathbb{R}^n . Formally,

$$M \text{ positive semi-definite} \iff \mathbf{x}^T M \mathbf{x} \geq 0 \text{ for all } \mathbf{x} \in \mathbb{R}^n$$

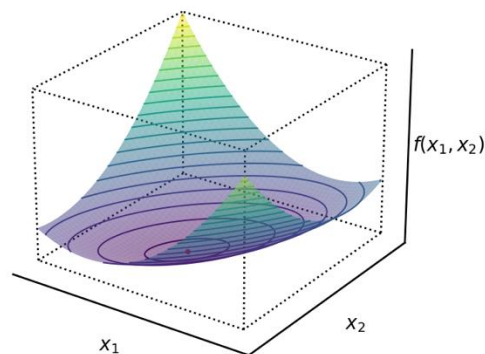
https://en.wikipedia.org/wiki/Definite_matrix

- Covariance matrix is positive semi-definite and symmetric.
 - Semi-definite because the variance can be zero

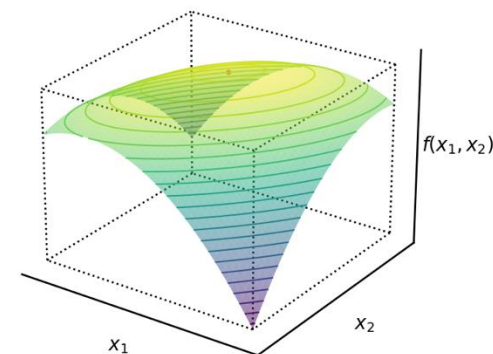
Positive semi-definite

$$f(x_1, x_2) = \frac{1}{2} \begin{bmatrix} x_1 & x_2 \end{bmatrix} \mathbf{A} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \mathbf{b}^\top \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + c$$

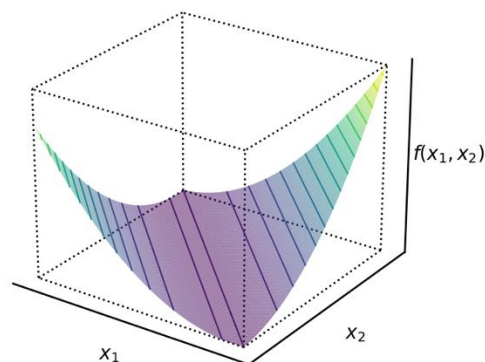
Positive definite



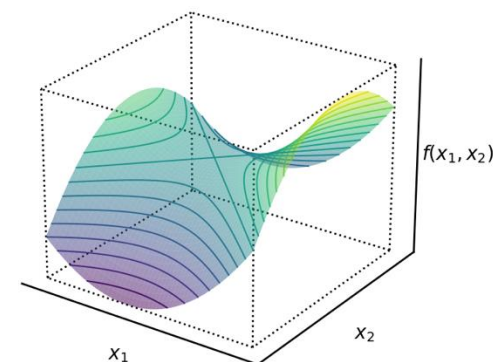
Negative definite



Singular & positive semi-definite



Indefinite



Positive semi-definite and covariance matrixes

Definitions for real matrices [\[edit \]](#)

An $n \times n$ symmetric real matrix M is said to be **positive-definite** if $\mathbf{x}^T M \mathbf{x} > 0$ for all non-zero \mathbf{x} in \mathbb{R}^n . Formally,

$$M \text{ positive-definite} \iff \mathbf{x}^T M \mathbf{x} > 0 \text{ for all } \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$$

An $n \times n$ symmetric real matrix M is said to be **positive-semidefinite** or **non-negative-definite** if $\mathbf{x}^T M \mathbf{x} \geq 0$ for all \mathbf{x} in \mathbb{R}^n . Formally,

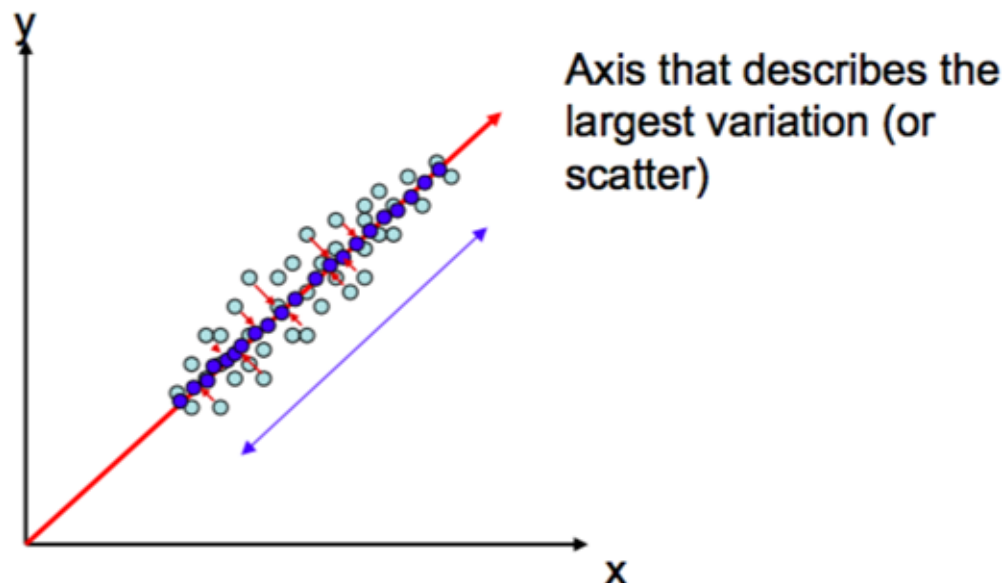
$$M \text{ positive semi-definite} \iff \mathbf{x}^T M \mathbf{x} \geq 0 \text{ for all } \mathbf{x} \in \mathbb{R}^n$$

https://en.wikipedia.org/wiki/Definite_matrix

- Covariance matrix is positive semi-definite and symmetric.
 - Symmetric \rightarrow real eigen values, eigen vectors are mutually orthogonal
 - $\mathbf{q}_i^T \mathbf{q}_j = 0$ for $i \neq j$
 - Positive semi-definite \rightarrow eigen values are nonnegative
 - Positive definite \rightarrow eigen values are positive \rightarrow invertible

What is PCA?

- We want to reduce the dimensionality but keep useful information
 - What is useful information? Variation
- We want to find a projection (a transformation) that describe maximum variation



Trace properties

$$1 \quad \bullet \operatorname{tr}(a) = a$$

$$2 \quad \bullet \operatorname{tr} A = \operatorname{tr} A^T$$

$$3 \quad \bullet \operatorname{tr}(A+B) = \operatorname{tr} A + \operatorname{tr} B$$

$$4 \quad \bullet \operatorname{tr}(aA) = a \operatorname{tr}(A)$$

$$5 \quad \nabla_A \operatorname{tr} AB = B^T$$

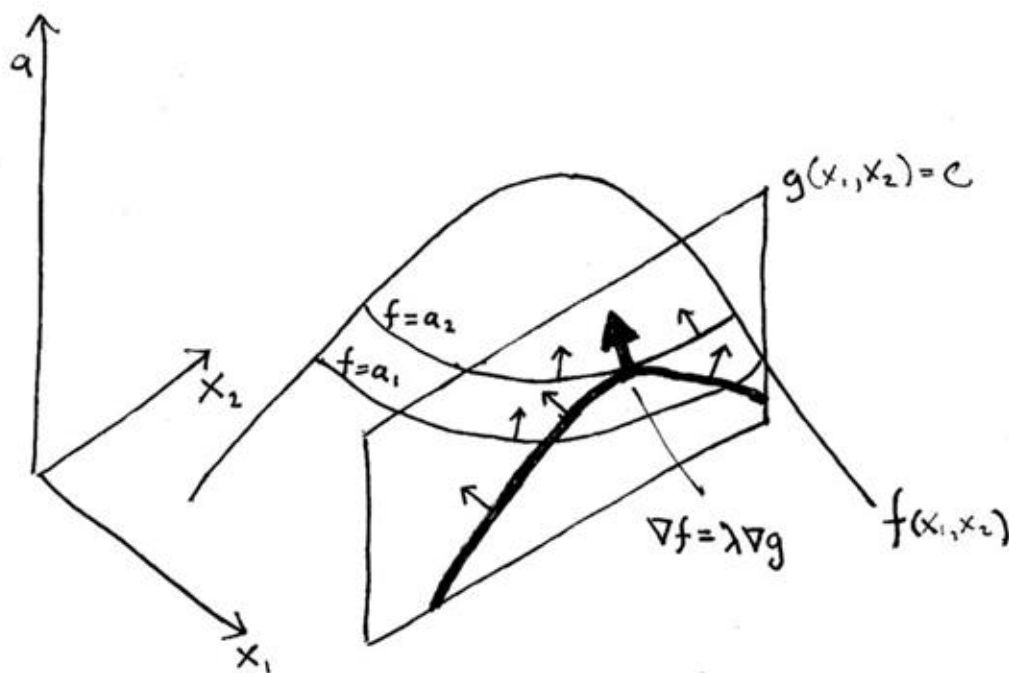
$$6 \quad \nabla_{A^T} f(A) = (\nabla_A f(A))^T$$

$$7 \quad \nabla_A \operatorname{tr} ABA^T C = CAB + C^T AB^T$$

$$8 \quad \nabla_{A^T} \operatorname{tr} ABA^T C = B^T A^T C^T + BA^T C$$

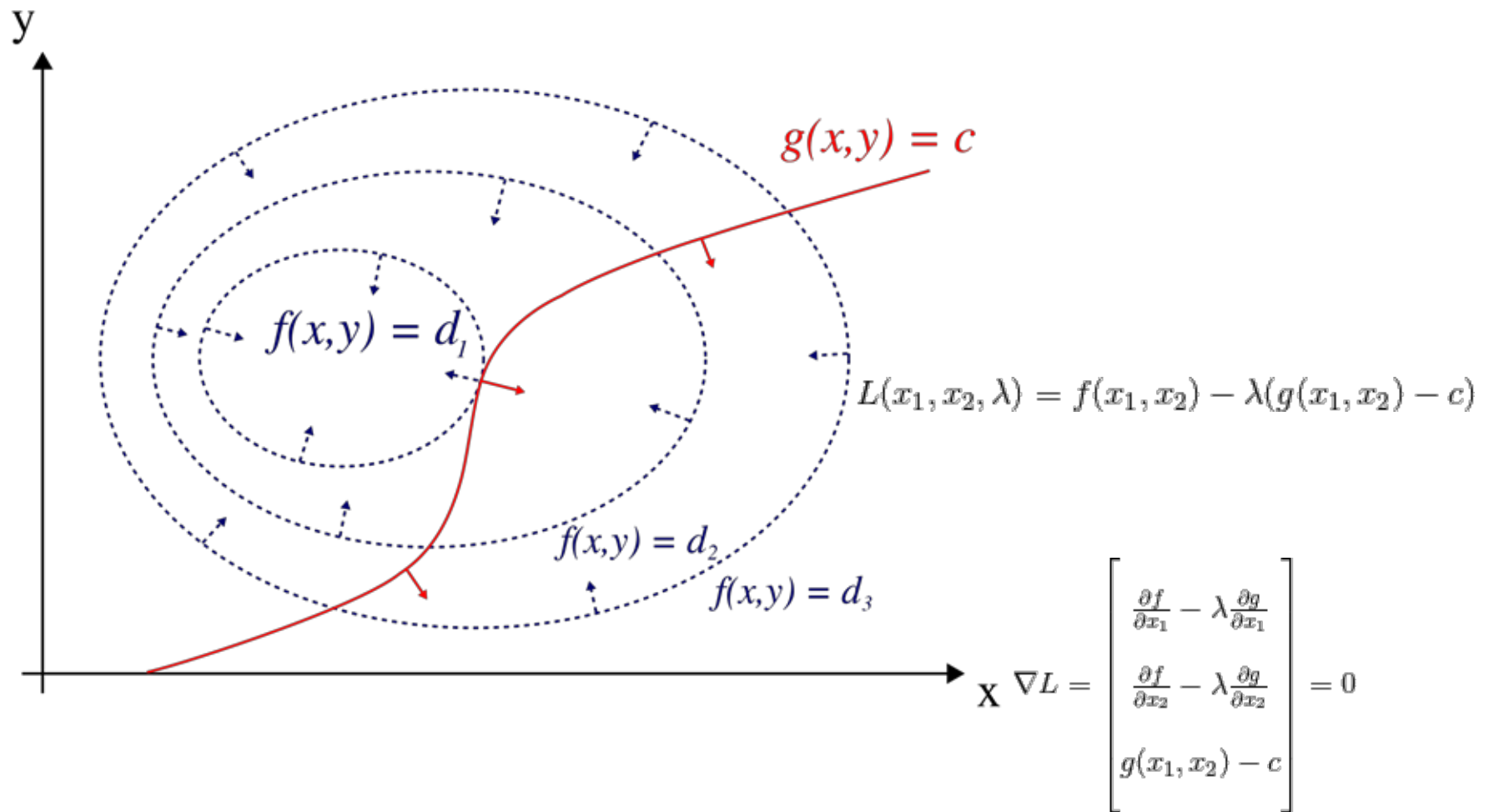
Lagrange Multiplier

<https://medium.com/@andrew.chamberlain/a-simple-explanation-of-why-lagrange-multipliers-works-253e2cdcbf74>



$$L(x_1, x_2, \lambda) = f(x_1, x_2) - \lambda(g(x_1, x_2) - c)$$

$$\nabla L = \begin{bmatrix} \frac{\partial f}{\partial x_1} - \lambda \frac{\partial g}{\partial x_1} \\ \frac{\partial f}{\partial x_2} - \lambda \frac{\partial g}{\partial x_2} \\ g(x_1, x_2) - c \end{bmatrix} = 0$$



So we got to eigenvectors

- A $d \times d$ covariance matrix has d eigenvectors/values pair. Do we use all of them?
- Which pair to use?

Selecting eigenvectors

- Remember the variance of projected data is

$$\omega^T \Sigma \omega. \quad (1)$$

- And our solution yielded $\Sigma \omega = \lambda \omega \quad (2)$

- Plug (2) in (1) and we get

$$\begin{aligned} \text{projected variance} &= \omega^T \Sigma \omega = \omega^T \lambda \omega \\ &= \lambda \omega^T \omega \quad (\text{remember } \|\omega\|=1) \\ &= \lambda \end{aligned}$$

PCA

- The direction vector captures the variance corresponding to the eigenvalue
- So we want the higher eigenvalues
 - How many?

Matrix rank

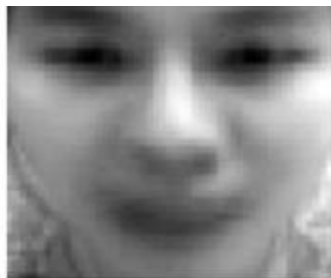
- A square $d \times d$ matrix has full rank (e.g. rank d) if the columns are linearly independent.
- The number of linearly independent columns is the rank of the matrix
- A covariance matrix of size $d \times d$ will have at most $N-1$ rank where N is the number of training samples
 - 640×640 images = ~ 400000 dimensions
 - 1000 training images
 - The covariance matrix will be at most rank 999. The missing rank is because of the mean.

PCA

- The direction vector captures the variance corresponding to the eigenvalue
- So we want the higher eigenvalues
- Take the eigenvalues with non-zero eigenvalues (at most $N-1$ non-zero eigenvalues)

Eigenfaces

Meanface



V1



V2



V3



V4



V5



V6



V7



V8



V9



V10

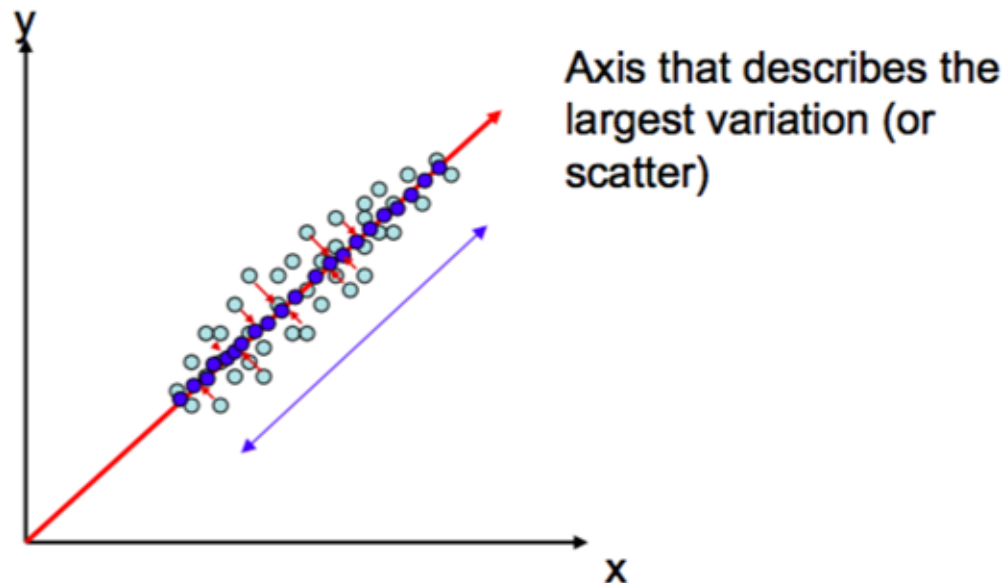


V11



What is PCA?

- We want to reduce the dimensionality but keep useful information
 - What is useful information? Variation
- We want to find a projection (a transformation) that describe maximum variation



Means

- In PCA, we model variance. (Variation around the mean)
- In our projection we need to remove the mean

$$\mathbf{p} = \mathbf{V}^T(\mathbf{x} - \mathbf{m})$$

- The mean is the mean of all your training data
- If we want to reconstruct the data we need to add back the mean

$$\mathbf{x} = \sum_{i=1}^N p_i \mathbf{v}_i + \mathbf{m} = p_1 \begin{bmatrix} | \\ \mathbf{v}_1 \\ | \end{bmatrix} + p_2 \begin{bmatrix} | \\ \mathbf{v}_2 \\ | \end{bmatrix} + \dots + p_n \begin{bmatrix} | \\ \mathbf{v}_n \\ | \end{bmatrix} + \mathbf{m} = \mathbf{V}\mathbf{p} + \mathbf{m}$$

Basis decomposition

- Let's consider our projection w_i which is the eigenvectors to be a **basis vector** v_i
- We can represent any vector as a sum of basis vectors as follows:

$$\mathbf{x} = \sum_{i=1}^N p_i \mathbf{v}_i = p_1 \begin{bmatrix} | \\ \mathbf{v}_1 \\ | \end{bmatrix} + p_2 \begin{bmatrix} | \\ \mathbf{v}_2 \\ | \end{bmatrix} + \dots + p_n \begin{bmatrix} | \\ \mathbf{v}_n \\ | \end{bmatrix} = \mathbf{V}\mathbf{p}$$

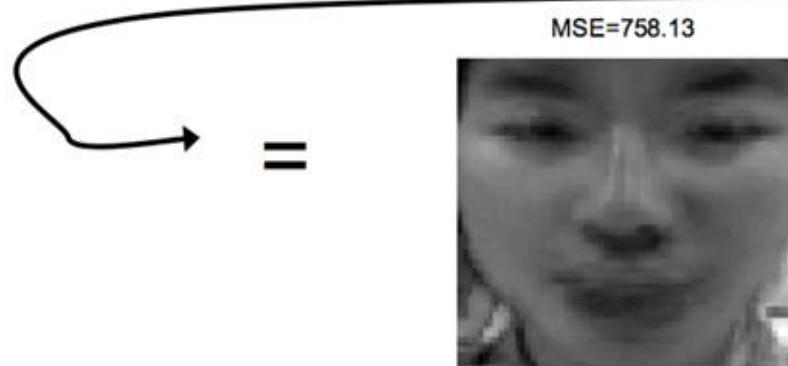
Finding the weights

$$\mathbf{x} = \sum_{i=1}^N p_i \mathbf{v}_i = p_1 \begin{bmatrix} | \\ \mathbf{v}_1 \\ | \end{bmatrix} + p_2 \begin{bmatrix} | \\ \mathbf{v}_2 \\ | \end{bmatrix} + \dots + p_n \begin{bmatrix} | \\ \mathbf{v}_n \\ | \end{bmatrix} = \mathbf{V} \mathbf{p}$$

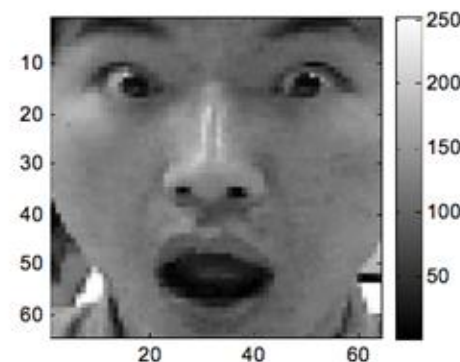
- If \mathbf{v}_i are orthogonal, the projection of \mathbf{x} onto \mathbf{v}_i gives p_i

$$\mathbf{V}^T \mathbf{x} = \begin{bmatrix} - & \mathbf{v}_1 & - \\ - & \mathbf{v}_2 & - \\ - & \mathbf{v}_3 & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

Reconstruction with eigenfaces



reconstructed
 $\tilde{\mathbf{X}}$



Original
 \mathbf{X}

Practical issues

- If your data has different magnitudes in different dimensions, normalize each dimension before PCA
- If we have 640x640 images = ~ 400000 dimensions.
- What is the size of the covariance matrix?
- 400000×400000



Practical issues

- You have N training examples.
- For the case where $N \ll 400000$, we only have $N-1$ eigenvalues we care about anyway

Gram Matrix

$$\Sigma = E(x - \mu)(x - \mu)^T = XX^T$$

Covariance matrix
is the **outer-product**
of the input matrix

Must solve $\Sigma v = \lambda v$

$$XX^T v = \lambda v \quad (\text{pre-mult by } X^T) \quad (1)$$

$$X^T X X^T v = \lambda X^T v \quad (v' = X^T v) \quad (2)$$

Solve eigenvalue problem $X^T X v' = \lambda v'$

- $X^T X$ is a gram of **inner-product** matrix. Its size is $N \times N$ where N is the number of data samples.

But how to get v from v' ?

- From previous slide, equation (1) and (2)
 - $XX^T v = \lambda v$ (1)
 - $v' = X^T v$ (2)
- Substitute (2) into (1)
 - $Xv' = \lambda v$
- Thus, $v = Xv'$. We don't care about the scaling term because we will always scale the eigenvector so that it is orthonormal i.e. $\|v\| = 1$.

How many eigenvectors?

- Select based on amount of variance explained
 - Sum of eigenvalues exceeds some percent of total
- Reconstruction error

$$\mathbf{x} = \sum_{i=1}^N p_i \mathbf{v}_i = p_1 \begin{bmatrix} | \\ \mathbf{v}_1 \\ | \end{bmatrix} + p_2 \begin{bmatrix} | \\ \mathbf{v}_2 \\ | \end{bmatrix} + \dots + p_n \begin{bmatrix} | \\ \mathbf{v}_n \\ | \end{bmatrix} = \mathbf{V}\mathbf{p}$$

- Select enough \mathbf{v} so that the difference between original \mathbf{x} and reconstructed \mathbf{x} is small

Summary

- Feature selection and dimensionality reduction techniques help combat the curse of dimensionality.
- PCA tries to maximize reconstruction of the data while using the least amount of dimension.