

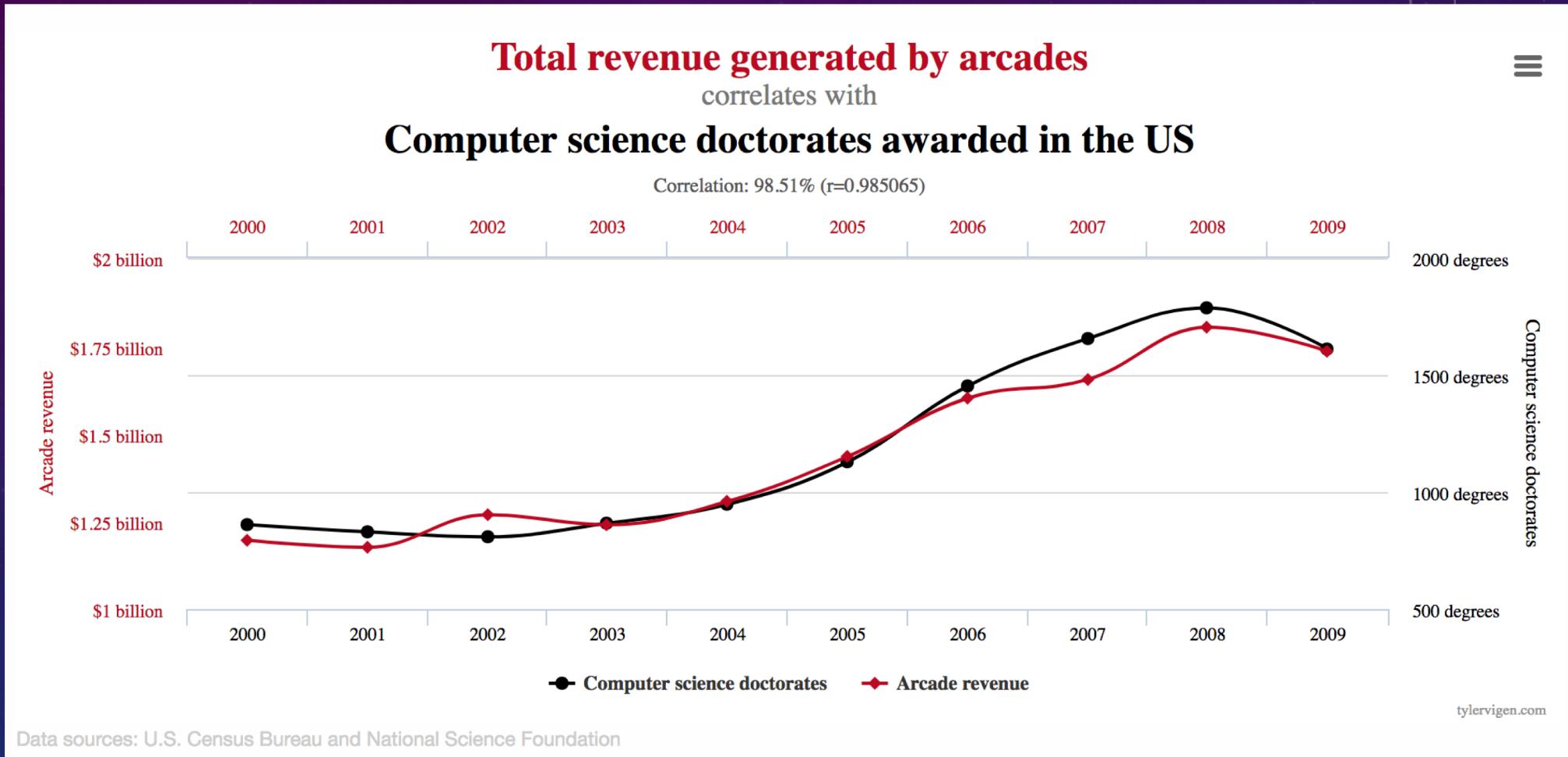
MATH FOUNDATIONS

8 JULY 2019 - WWW.SEA-MLS.COM

CHENG SOON ONG

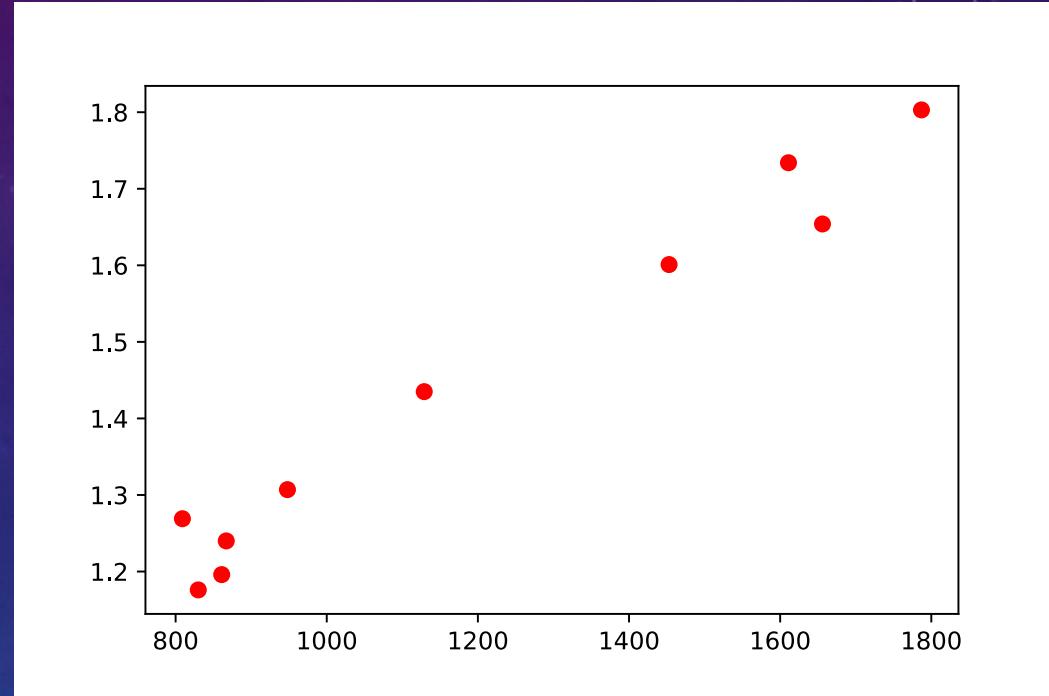
ong-home.my

PREDICT ARCADE REVENUE FROM CS PHD



PREDICT ARCADE REVENUE FROM CS PHD

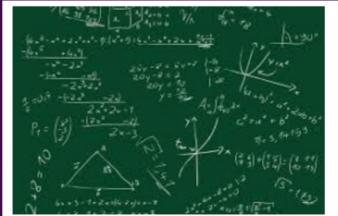
Year	CS PhD in USA	Arcade (\$billions)
2000	861	1.196
2001	830	1.176
2002	809	1.269
2003	867	1.240
2004	948	1.307
2005	1129	1.435
2006	1453	1.601
2007	1656	1.654
2008	1787	1.803
2009	1611	1.734



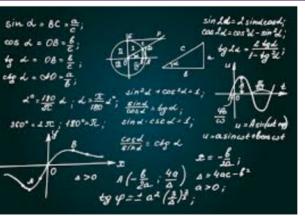
tylervigen.com

Looks like machine learning, but before we can learn that we need to do some maths

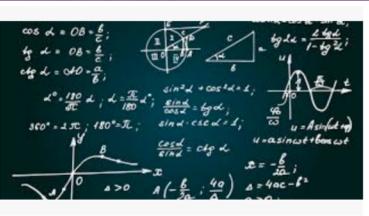
MATHEMATICS? THIS IS BORING...



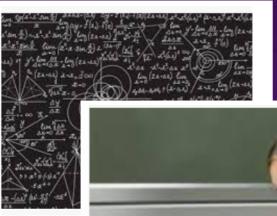
What is Mathematics?
livescience.com



Mathematics: forget simplicity, the ...
theconversation.com



Mathematics: forget simplicity, the ...
theconversation.com



Yes, mathem...
theconversat...



Students Sleeping in Class ...
blogs.edweek.org



Techniques for Sleeping in Class – The ...
sites.imsa.edu



Grades Suffer When Class Time Doesn't ...
blogs.edweek.org



Master's programme in Mathematics | KTH ...
kth.se



M.Sc. Mathematics - Minhaj Universi...
mul.edu.pk



Pathways: Mathematics | Yeshiva University
yu.edu



Mathematics - The Languag...
medium.com



Professor gives extra credit to ...
washingtonexaminer.com



Avoid Getting Caught Slee...
boredpanda.com



Wake Up to a Back-to-School Sleep ...
blog.nemours.org



Avoid Ge...
boredpan...



The Cure for Students Sleeping in Class



Sleeping Through the Semester: A St...
Sleep...



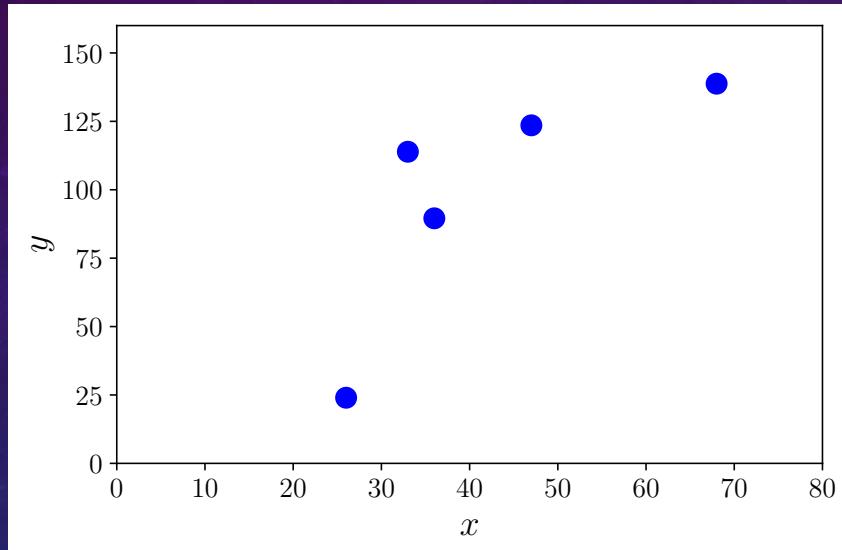
Want to take a nap after class ...
Sleep...



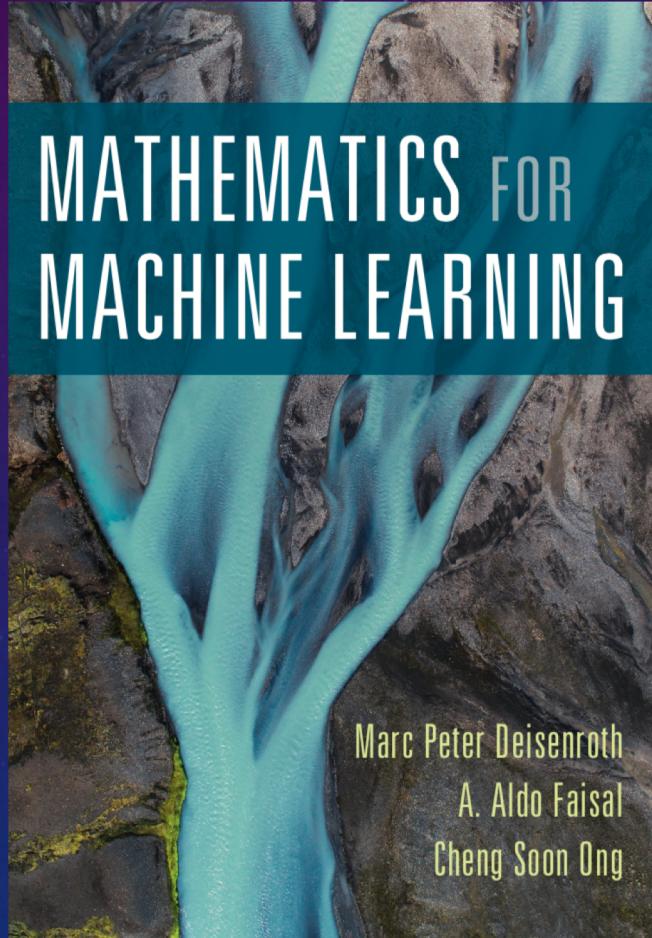
Sleep...

SHAMELESS PLUG: MATH FOR ML

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



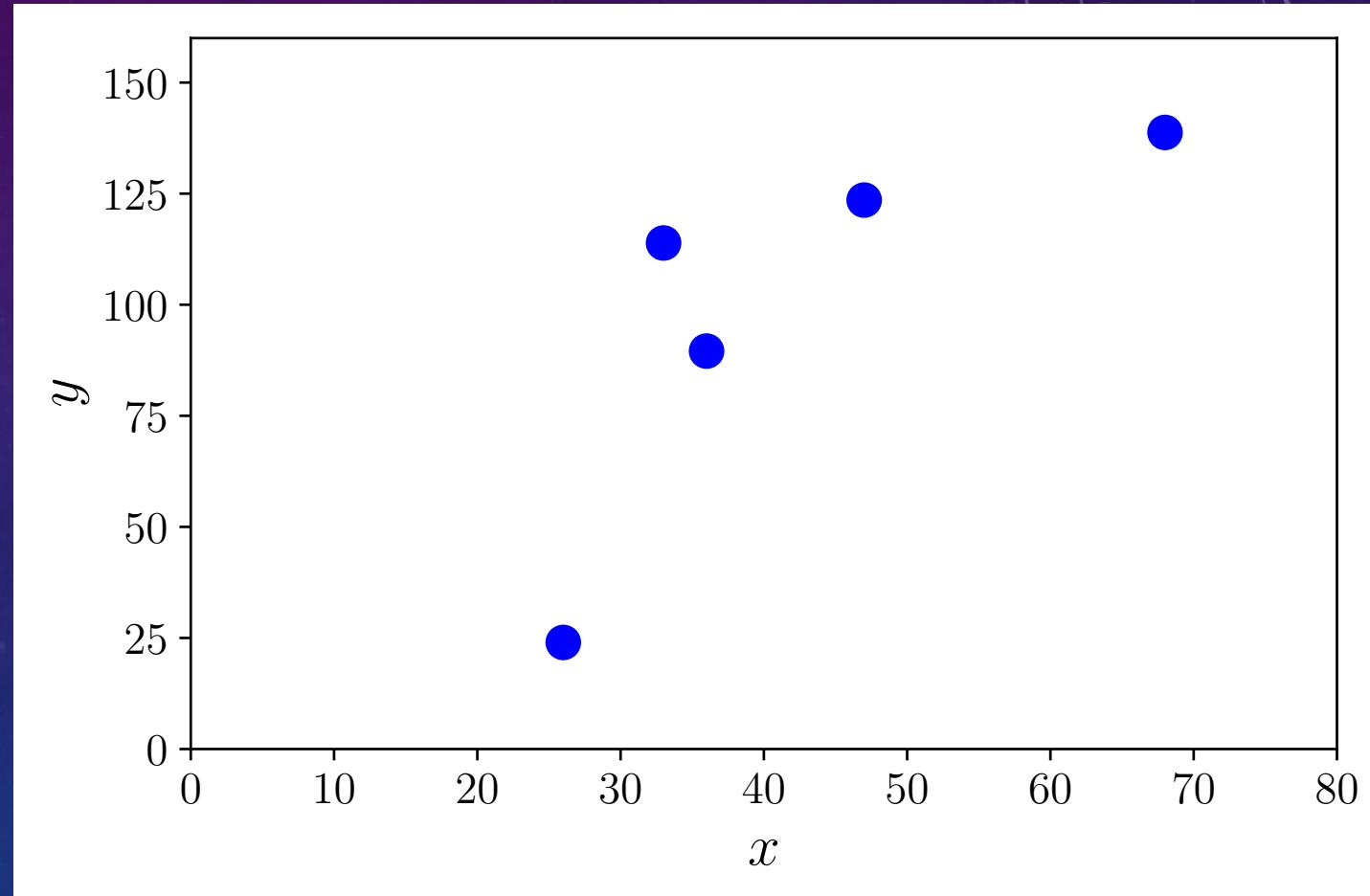
Predict salary (y) from age (x)



WHAT DO WE MEAN WHEN WE DRAW DOTS?

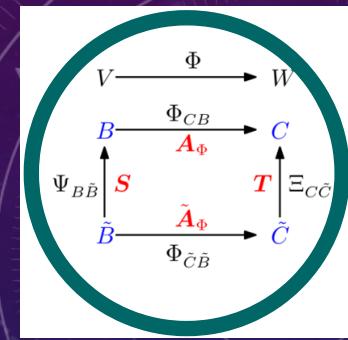
- Three views of a vector
 - (CS) array of numbers
 - (physics) magnitude and direction
 - (math) satisfies + and \times

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



LINEAR ALGEBRA – WHAT IS A VECTOR?

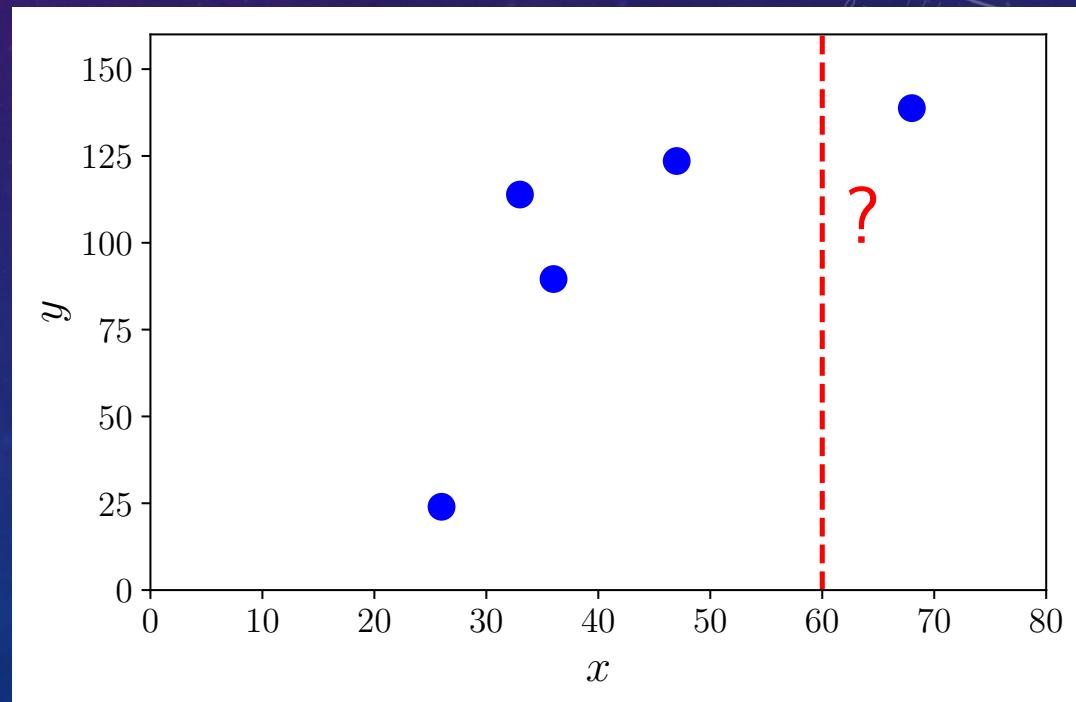
- Algebra: Set of objects and set of rules to manipulate them
 - Objects: vectors x and y
 - Rules: $+$ and \times , as well as defining a zero.
- Linear: $ax + by$
 - distributivity
 - associativity
- Vector space:
 - Closure: adding and scaling vectors keeps things in the vector space



MACHINE LEARNING IS ABOUT PREDICTION

- The values of y for the training data is not the main focus
- We are interested in generalization error:
 - What is the error we make on unseen data?
- Do not train on the test set

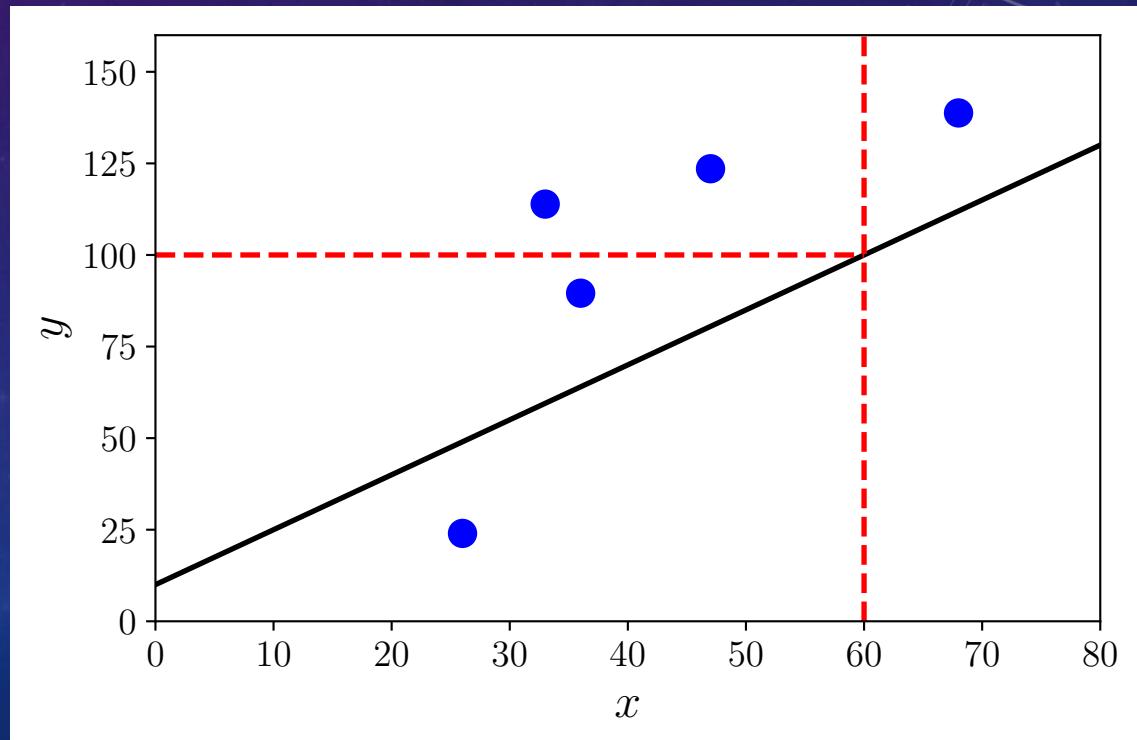
Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



FITTING A LINE

- The values of y for the training data is not the main focus
- We are interested in generalization error:
 - What is the error we make on unseen data?
- Do not train on the test set

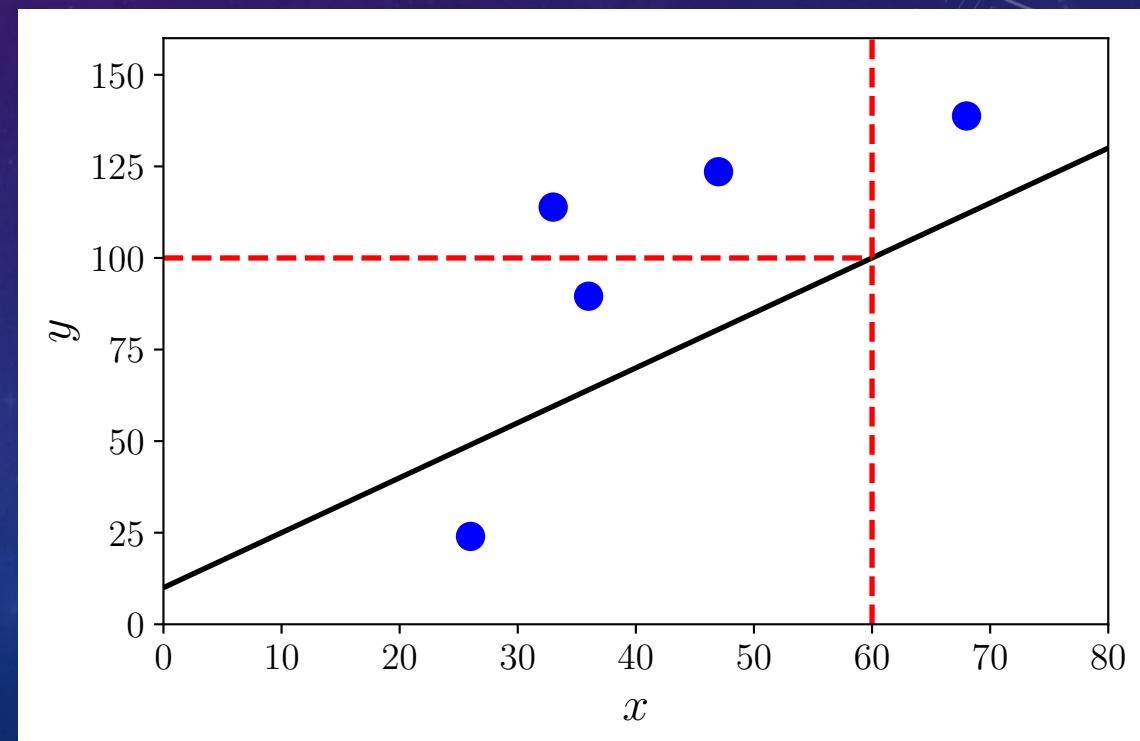
Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



FITTING A LINE - NOTATION

- $N = 5$
- $\mathbf{x} = [x_1, \dots, x_N]^T$
- $\mathbf{y} = [y_1, \dots, y_N]^T$
- x_n is a real number
- y_n is a real number
- $f(x) = w_1 x + w_0$
- $\mathbf{X} = [\mathbf{x} \ 1]$
- Find the best line that fits the data

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



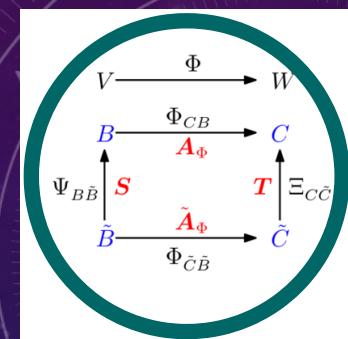
LINEAR ALGEBRA - MATRIX

- Want to solve $Xw = y$
- Solutions of linear equations
- Inverse and transpose

Definition 2.3 (Inverse). Consider a square matrix $A \in \mathbb{R}^{n \times n}$. Let matrix $B \in \mathbb{R}^{n \times n}$ have the property that $AB = I_n = BA$. B is called the *inverse* of A and denoted by A^{-1} .

Definition 2.4 (Transpose). For $A \in \mathbb{R}^{m \times n}$ the matrix $B \in \mathbb{R}^{n \times m}$ with $b_{ij} = a_{ji}$ is called the *transpose* of A . We write $B = A^\top$.

$$\begin{aligned}AA^{-1} &= I = A^{-1}A \\(AB)^{-1} &= B^{-1}A^{-1} \\(A + B)^{-1} &\neq A^{-1} + B^{-1} \\(A^\top)^\top &= A \\(A + B)^\top &= A^\top + B^\top \\(AB)^\top &= B^\top A^\top\end{aligned}$$



LINEAR ALGEBRA

- Linear independence

Definition 2.11 (Linear Combination). Consider a vector space V and a finite number of vectors $\mathbf{x}_1, \dots, \mathbf{x}_k \in V$. Then, every $\mathbf{v} \in V$ of the form

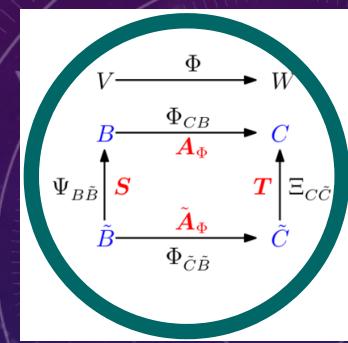
$$\mathbf{v} = \lambda_1 \mathbf{x}_1 + \cdots + \lambda_k \mathbf{x}_k = \sum_{i=1}^k \lambda_i \mathbf{x}_i \in V \quad (2.65)$$

with $\lambda_1, \dots, \lambda_k \in \mathbb{R}$ is a *linear combination* of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_k$.

- Basis and rank

Definition 2.14 (Basis). Consider a vector space $V = (\mathcal{V}, +, \cdot)$ and $\mathcal{A} \subseteq \mathcal{V}$. A generating set \mathcal{A} of V is called *minimal* if there exists no smaller set $\tilde{\mathcal{A}} \subseteq \mathcal{A} \subseteq \mathcal{V}$ that spans V . Every linearly independent generating set of V is minimal and is called a *basis* of V .

- Matrix: represent data vs represent transformations
- Linear vs Affine space: what is a linear regressor?

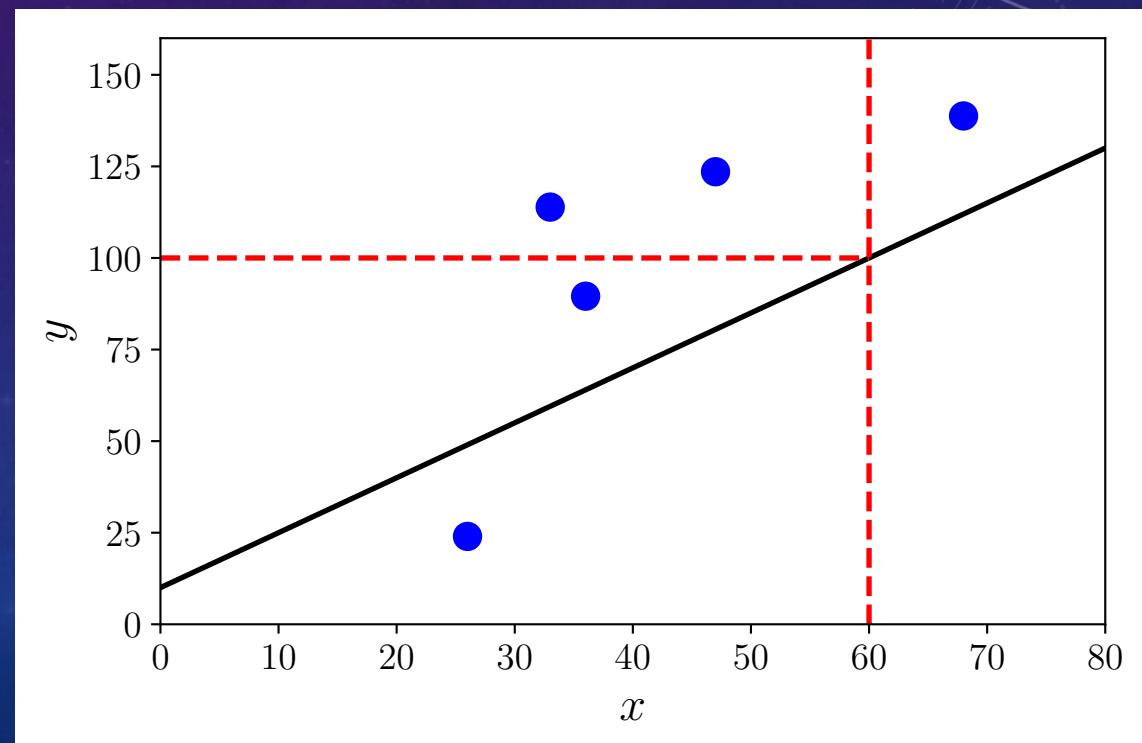


FITTING A LINE – LINEAR ALGEBRA

- Want to solve $Xw = y$
- If points don't fall perfectly on the line, no solution
- Find a point z that lies in the column space of X and is closest to y

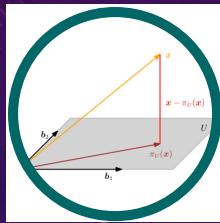
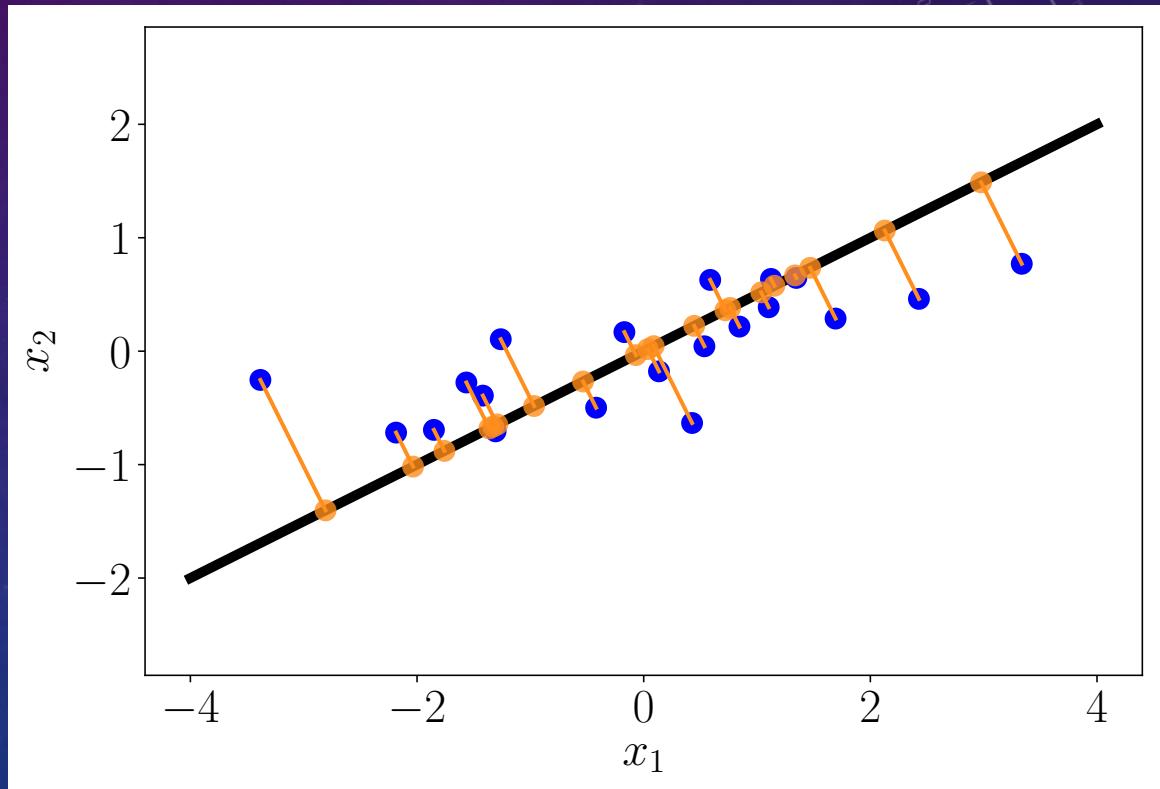
Closest?

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



ANALYTIC GEOMETRY

- Inner products
- Distances
- Orthogonality
- Orthogonal projection



ANALYTIC GEOMETRY

- Inner products

Definition 3.2. Let V be a vector space and $\Omega : V \times V \rightarrow \mathbb{R}$ be a bilinear mapping that takes two vectors and maps them onto a real number. Then

- Ω is called *symmetric* if $\Omega(\mathbf{x}, \mathbf{y}) = \Omega(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in V$, i.e., the order of the arguments does not matter.
- Ω is called *positive definite* if

$$\forall \mathbf{x} \in V \setminus \{\mathbf{0}\} : \Omega(\mathbf{x}, \mathbf{x}) > 0, \quad \Omega(\mathbf{0}, \mathbf{0}) = 0. \quad (3.8)$$

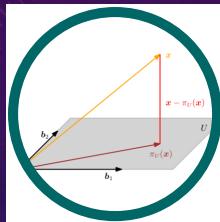
A positive definite, symmetric bilinear mapping $\Omega : V \times V \rightarrow \mathbb{R}$ is called an *inner product* on V . We typically write $\langle \mathbf{x}, \mathbf{y} \rangle$ instead of $\Omega(\mathbf{x}, \mathbf{y})$.

- Distances

$$\|\mathbf{x}\| := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$$

Definition 3.6 (Distance and Metric). Consider an inner product space $(V, \langle \cdot, \cdot \rangle)$. Then

$$d(\mathbf{x}, \mathbf{y}) := \|\mathbf{x} - \mathbf{y}\| = \sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle} \quad (3.21)$$



ANALYTIC GEOMETRY

- Orthogonality

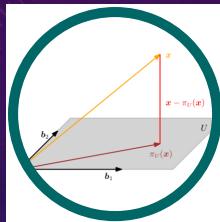
Definition 3.7 (Orthogonality). Two vectors \mathbf{x} and \mathbf{y} are *orthogonal* if and only if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$, and we write $\mathbf{x} \perp \mathbf{y}$. If additionally $\|\mathbf{x}\| = 1 = \|\mathbf{y}\|$, i.e., the vectors are unit vectors, then \mathbf{x} and \mathbf{y} are *orthonormal*.

- Orthogonal projection (recall linear mapping = transformation matrix)

Definition 3.10 (Projection). Let V be a vector space and $U \subseteq V$ a subspace of V . A linear mapping $\pi : V \rightarrow U$ is called a *projection* if $\pi^2 = \pi \circ \pi = \pi$.

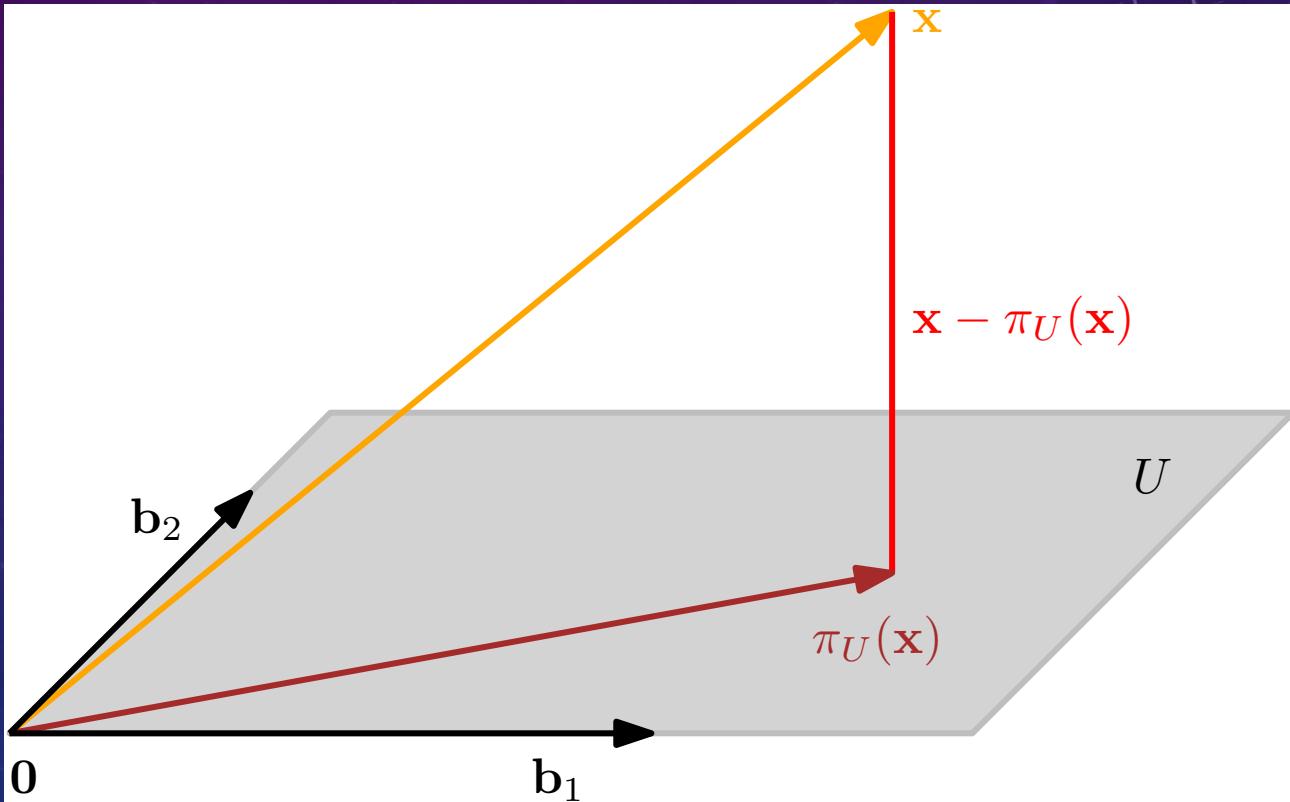
The projection $\pi_U(\mathbf{x})$ is closest to \mathbf{x} , where “closest” implies that the distance $\|\mathbf{x} - \pi_U(\mathbf{x})\|$ is minimal. It follows that the segment $\pi_U(\mathbf{x}) - \mathbf{x}$ from $\pi_U(\mathbf{x})$ to \mathbf{x} is orthogonal to U , and therefore the basis vector \mathbf{b} of U . The orthogonality condition yields $\langle \pi_U(\mathbf{x}) - \mathbf{x}, \mathbf{b} \rangle = 0$ since angles between vectors are defined via the inner product.

The projection $\pi_U(\mathbf{x})$ of \mathbf{x} onto U must be an element of U and, therefore, a multiple of the basis vector \mathbf{b} that spans U . Hence, $\pi_U(\mathbf{x}) = \lambda \mathbf{b}$, for some $\lambda \in \mathbb{R}$.



ANALYTIC GEOMETRY

- Want to solve $Xw = y$
- Find a point z that lies in the column space of X and is closest to y
- z is found by the orthogonal projection of y onto the column space of X



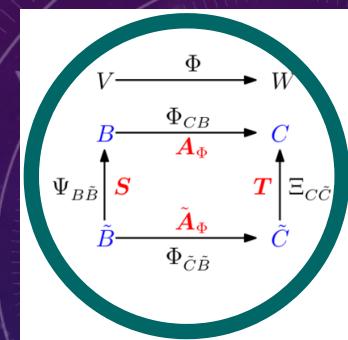
REMINDER OF MATRIX OPERATIONS

- Want to solve $\mathbf{X}\mathbf{w} = \mathbf{y}$
- Solutions of linear equations
- Inverse and transpose

Definition 2.3 (Inverse). Consider a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$. Let matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ have the property that $\mathbf{AB} = \mathbf{I}_n = \mathbf{BA}$. \mathbf{B} is called the *inverse* of \mathbf{A} and denoted by \mathbf{A}^{-1} .

Definition 2.4 (Transpose). For $\mathbf{A} \in \mathbb{R}^{m \times n}$ the matrix $\mathbf{B} \in \mathbb{R}^{n \times m}$ with $b_{ij} = a_{ji}$ is called the *transpose* of \mathbf{A} . We write $\mathbf{B} = \mathbf{A}^\top$.

$$\begin{aligned} \mathbf{AA}^{-1} &= \mathbf{I} = \mathbf{A}^{-1}\mathbf{A} \\ (\mathbf{AB})^{-1} &= \mathbf{B}^{-1}\mathbf{A}^{-1} \\ (\mathbf{A} + \mathbf{B})^{-1} &\neq \mathbf{A}^{-1} + \mathbf{B}^{-1} \\ (\mathbf{A}^\top)^\top &= \mathbf{A} \\ (\mathbf{A} + \mathbf{B})^\top &= \mathbf{A}^\top + \mathbf{B}^\top \\ (\mathbf{AB})^\top &= \mathbf{B}^\top \mathbf{A}^\top \end{aligned}$$

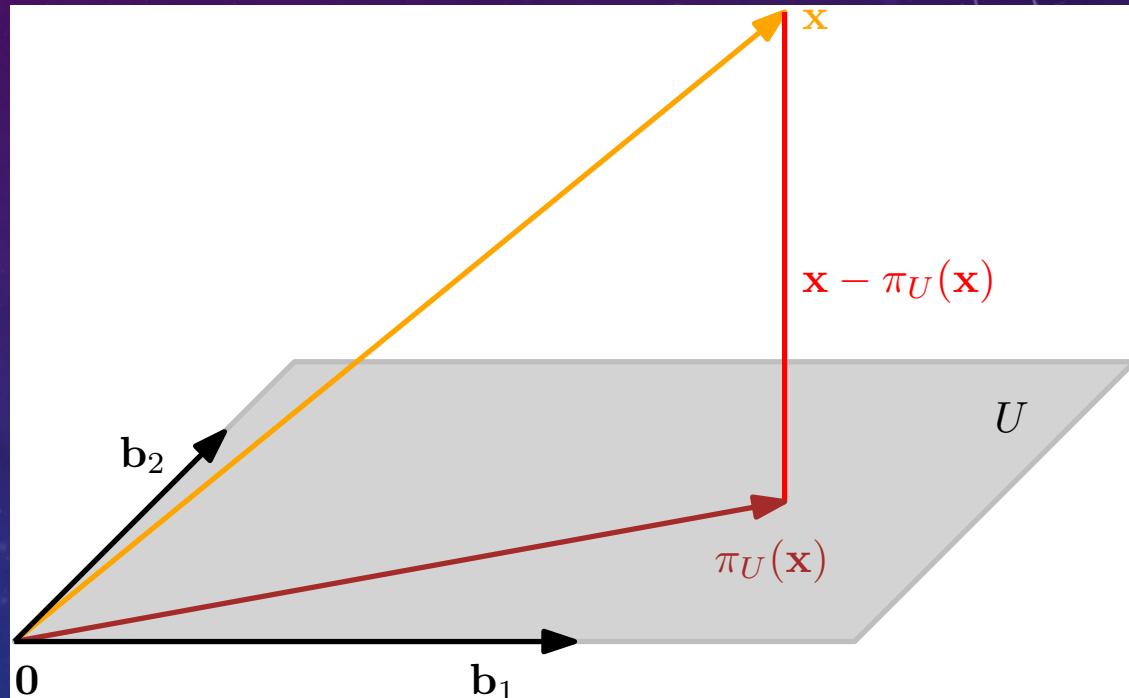


ANALYTIC GEOMETRY

- z is found by the orthogonal projection of y onto the column space of X
- Column space of X is spanned by $\{x, 1\}$, and hence we need to find coordinates w_1 and w_0 of the projection, such that the linear combination Xw is closest to y .
- Closest means that the vector connecting z to y is orthogonal to the column space of X .

$$X^T(y - z) = 0 \rightarrow X^T(y - Xw) = 0$$

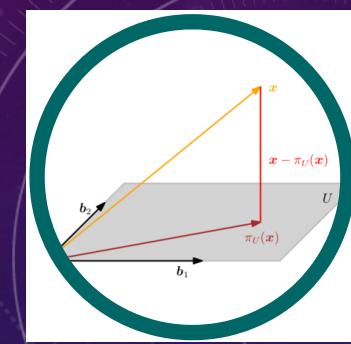
- Solving gives $w = (X^T X)^{-1} X^T y$



ANALYTIC GEOMETRY

- Want to solve $Xw = y$
- Find a point z that lies in the column space of X and is closest to y
- z is found by the orthogonal projection of y onto the column space of X

- $N = 5$
- $x = [x_1, \dots, x_N]^T$
- $y = [y_1, \dots, y_N]^T$
- x_n is a real number
- y_n is a real number
- $f(x) = w_1 x + w_0$
- $X = [x \ 1]$
- $w^* = (X^T X)^{-1} X^T y$



MATRIX DECOMPOSITIONS

- How do we compute the inverse of a matrix?

$$\mathbf{A}^{-1} = \frac{1}{a_{11}a_{22} - a_{12}a_{21}} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}$$

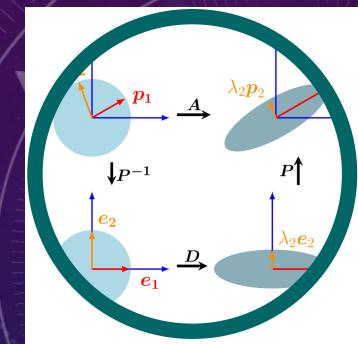
$$\det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}$$

- But for matrices that are larger, we do not have a closed form rule.
- Recall that linear mappings have an associated transformation matrix
- Disentangle different parts by an eigenvalue decomposition (inverse of a diagonal matrix is easy)

Theorem 4.20 (Eigendecomposition). A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be factored into

$$\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}, \quad (4.55)$$

where $\mathbf{P} \in \mathbb{R}^{n \times n}$ and \mathbf{D} is a diagonal matrix whose diagonal entries are the eigenvalues of \mathbf{A} , if and only if the eigenvectors of \mathbf{A} form a basis of \mathbb{R}^n .



According to the Abel–Ruffini theorem, there is in general no algebraic solution for polynomials of degree 5 or more (Abel, 1826).

OTHER MATRIX DECOMPOSITIONS

- For positive definite matrices

$$\forall \mathbf{x} \in V \setminus \{\mathbf{0}\} : \mathbf{x}^\top \mathbf{A} \mathbf{x} > 0.$$

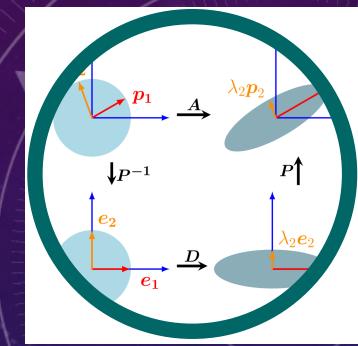
we have the Cholesky decomposition

Theorem 4.18 (Cholesky Decomposition). *A symmetric, positive definite matrix \mathbf{A} can be factorized into a product $\mathbf{A} = \mathbf{L}\mathbf{L}^\top$, where \mathbf{L} is a lower-triangular matrix with positive diagonal elements:*

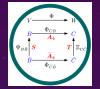
- For non-square matrices we have the singular value decomposition

Theorem 4.22 (SVD Theorem). *Let $\mathbf{A}^{m \times n}$ be a rectangular matrix of rank $r \in [0, \min(m, n)]$. The SVD of \mathbf{A} is a decomposition of the form*

$$_m^n \mathbf{A} = _m^m \mathbf{U} _m^n \Sigma _n^n \mathbf{V}^\top u \quad (4.64)$$

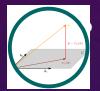


SUMMARY: MATH FOR ML



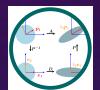
Chap. 2: Linear Algebra

vector space, linear maps, affine space



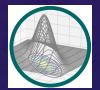
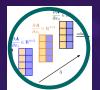
Chap. 3: Analytic Geometry

inner products, distances, orthogonality



Chap. 4: Matrix Decompositions

Cholesky, eigenvalue decomposition, singular value decomposition



MATHEMATICS FOR MACHINE LEARNING



mml-book.com

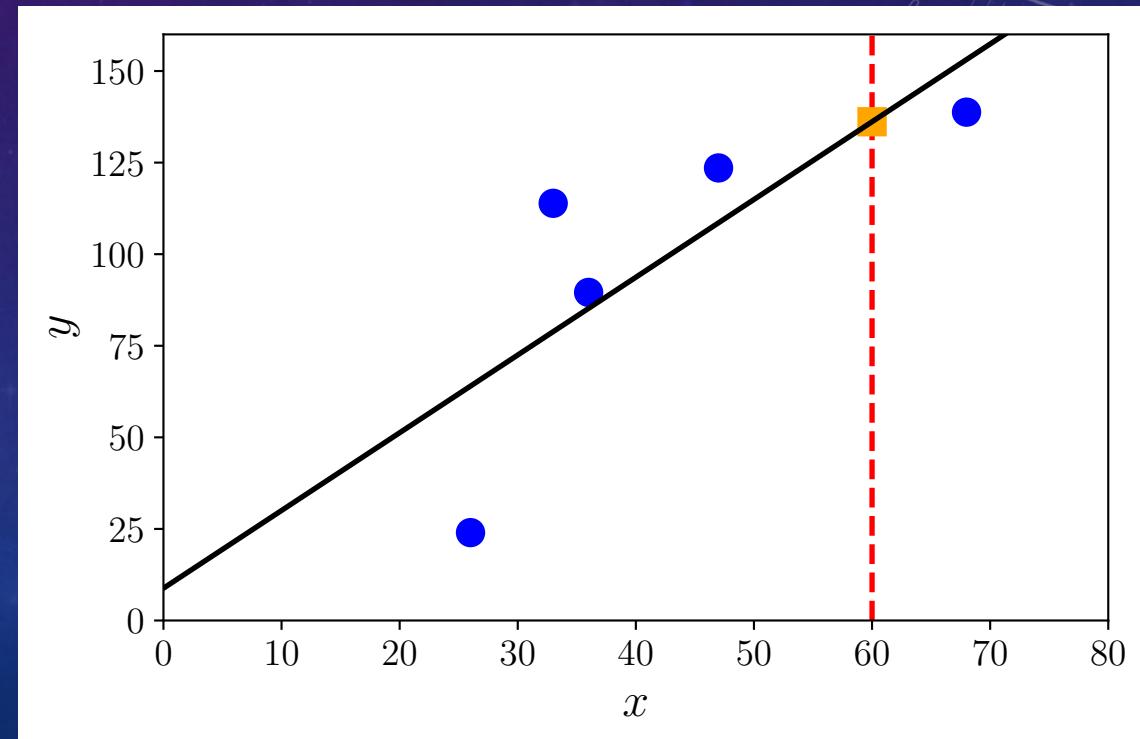
FITTING A LINE – OPTIMIZATION

- Want to solve $Xw = y$
- If points don't fall perfectly on the line, no solution
- Instead, find the closest approximate solution

$$\min_w || Xw - y ||^2$$

- Solve for a minimum by taking the gradient and setting to zero.

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



VECTOR CALCULUS - GRADIENT

$$\ell(x) = x^4 + 7x^3 + 5x^2 - 17x + 3,$$

Univariate calculus

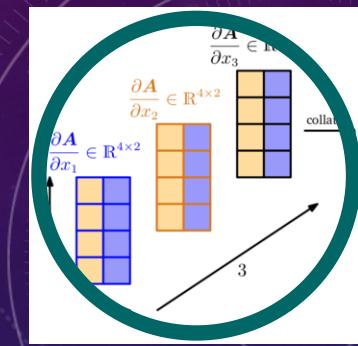
$$\frac{d\ell(x)}{dx} = 4x^3 + 21x^2 + 10x - 17.$$

Definition 5.5 (Partial Derivative). For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^n$ of n variables x_1, \dots, x_n we define the *partial derivatives* as

$$\begin{aligned} \frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h} \\ &\vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(\mathbf{x})}{h} \end{aligned} \tag{5.39}$$

and collect them in the row vector

$$\nabla_{\mathbf{x}} f = \text{grad } f = \frac{df}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}, \tag{5.40}$$



row vector

FITTING A LINE – OPTIMIZATION

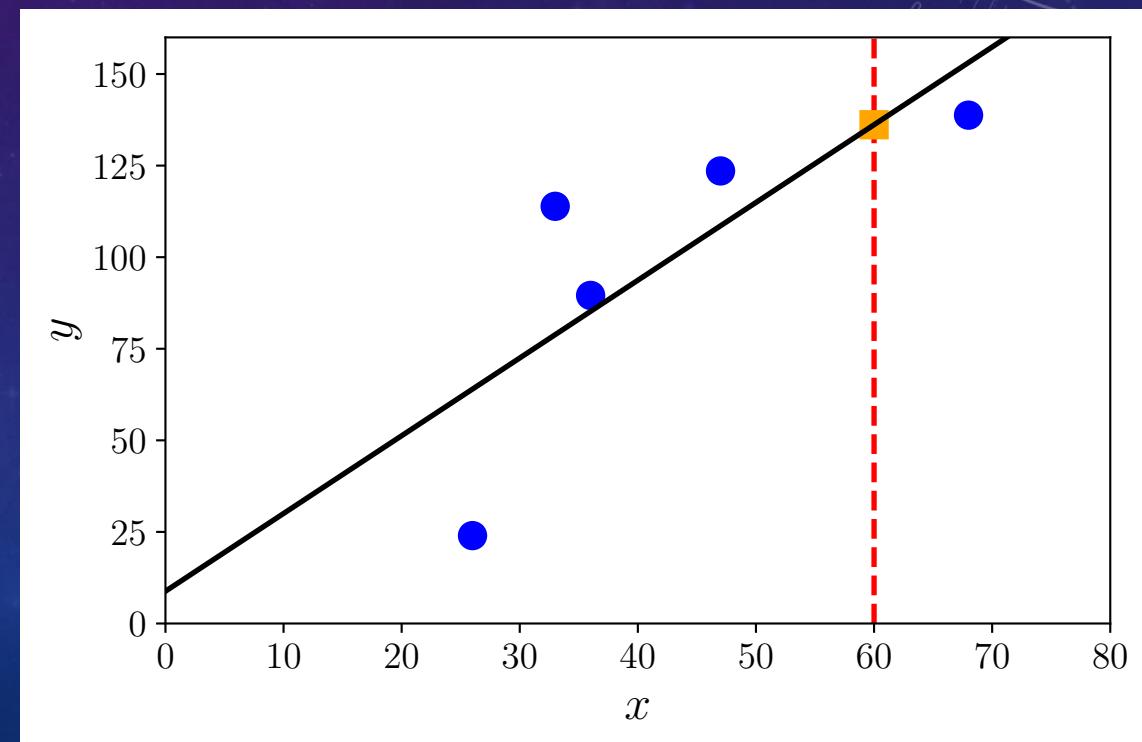
- Find the closest approximate solution

$$\min_w || Xw - y ||^2$$

- Solve for a minimum by taking the gradient and setting to zero.
- Gradient (wrt w) is $2 (Xw - y)^T X$
- Solving for stationary point gives

$$X^T X w = X^T y$$

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



VECTOR CALCULUS - JACOBIAN

Vector valued functions

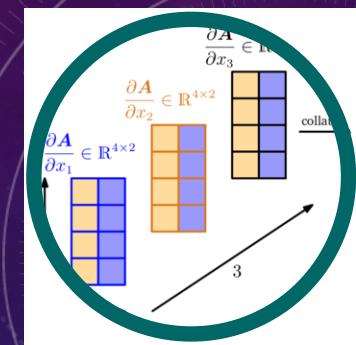
$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^m.$$

$$\frac{\partial \mathbf{f}}{\partial x_i} = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix}$$

Definition 5.6 (Jacobian). The collection of all first-order partial derivatives of a vector-valued function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called the *Jacobian*. The Jacobian \mathbf{J} is an $m \times n$ matrix, which we define and arrange as follows:

$$\mathbf{J} = \nabla_{\mathbf{x}} \mathbf{f} = \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \left[\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} \quad \dots \quad \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \right] \quad (5.57)$$

The gradient of a function
 $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a matrix of size $m \times n$.



SUM RULE, PRODUCT RULE, CHAIN RULE

Product rule:

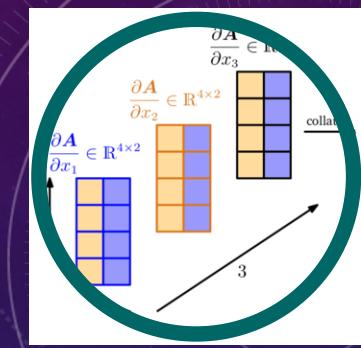
$$(fg)' = f'g + fg',$$

Sum rule:

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

Chain rule:

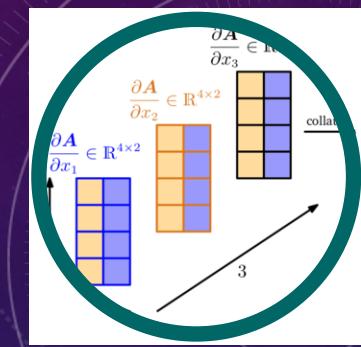
$$(g(f))' = g'(f)f'$$



SUM RULE, PRODUCT RULE, CHAIN RULE

Product rule: $\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g}{\partial \mathbf{x}}$ (5.46)

Sum rule: $\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$ (5.47)



Product rule:
 $(fg)' = f'g + fg'$,
Sum rule:
 $(f + g)' = f' + g'$,
Chain rule:
 $(g(f))' = g'(f)f'$

SUM RULE, PRODUCT RULE, CHAIN RULE

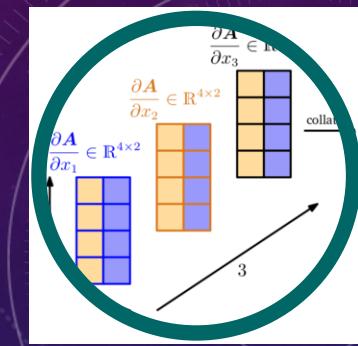
Product rule: $\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g}{\partial \mathbf{x}}$ (5.46)

Sum rule: $\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$ (5.47)

If $f(x_1, x_2)$ is a function of x_1 and x_2 , where $x_1(s, t)$ and $x_2(s, t)$ themselves functions of two variables s and t , the chain rule yields

$$\begin{aligned}\frac{\partial f}{\partial s} &= \frac{\partial f}{\partial \mathbf{x}_1} \frac{\partial \mathbf{x}_1}{\partial s} + \frac{\partial f}{\partial \mathbf{x}_2} \frac{\partial \mathbf{x}_2}{\partial s}, \\ \frac{\partial f}{\partial t} &= \frac{\partial f}{\partial \mathbf{x}_1} \frac{\partial \mathbf{x}_1}{\partial t} + \frac{\partial f}{\partial \mathbf{x}_2} \frac{\partial \mathbf{x}_2}{\partial t},\end{aligned}$$

$$\begin{aligned}\frac{df}{d(s, t)} &= \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)} = \underbrace{\left[\frac{\partial f}{\partial \mathbf{x}_1} \quad \frac{\partial f}{\partial \mathbf{x}_2} \right]}_{= \frac{\partial f}{\partial \mathbf{x}}} \underbrace{\left[\begin{array}{c|c} \frac{\partial \mathbf{x}_1}{\partial s} & \frac{\partial \mathbf{x}_1}{\partial t} \\ \hline \frac{\partial \mathbf{x}_2}{\partial s} & \frac{\partial \mathbf{x}_2}{\partial t} \end{array} \right]}_{= \frac{\partial \mathbf{x}}{\partial (s, t)}} \\ &= \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)}\end{aligned}$$

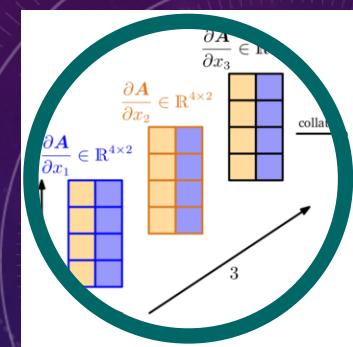


Product rule:
 $(fg)' = f'g + fg'$,
 Sum rule:
 $(f + g)' = f' + g'$,
 Chain rule:
 $(g(f))' = g'(f)f'$

CHAIN RULE

If $f(x_1, x_2)$ is a function of x_1 and x_2 , where $x_1(s, t)$ and $x_2(s, t)$ themselves functions of two variables s and t , the chain rule yields

$$\frac{df}{d(s, t)} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)} = \underbrace{\begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix}}_{\frac{\partial f}{\partial \mathbf{x}}} \underbrace{\begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{bmatrix}}_{\frac{\partial \mathbf{x}}{\partial (s, t)}}$$



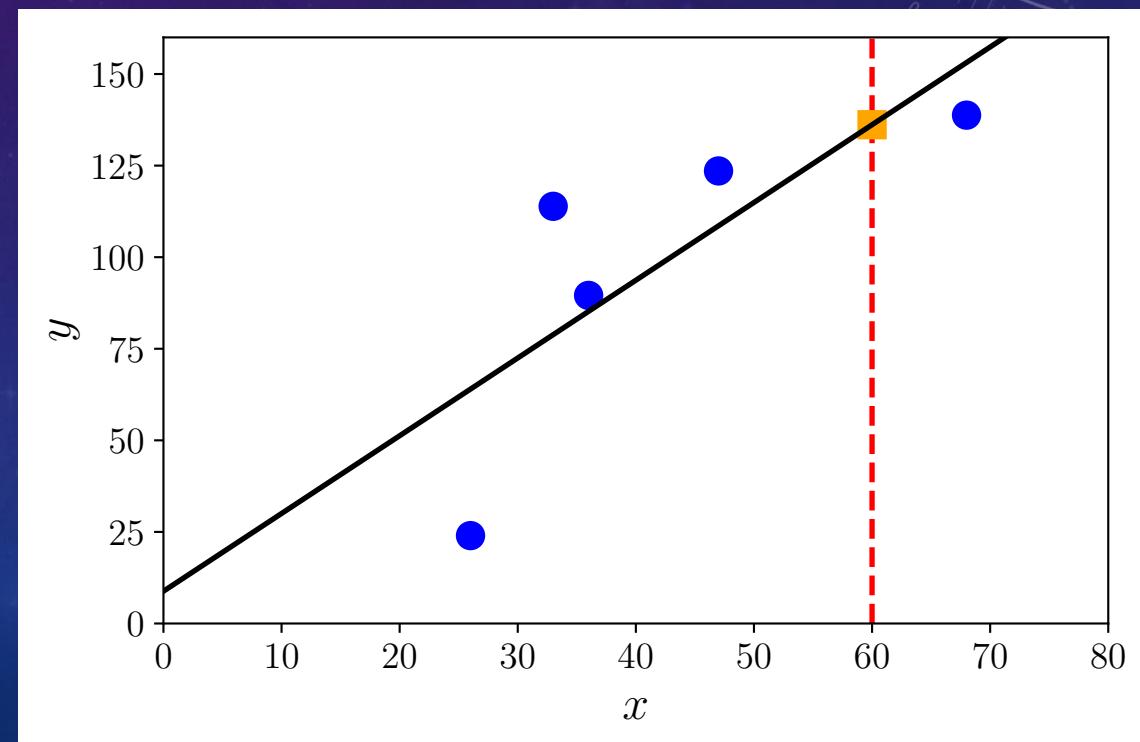
FITTING A LINE – OPTIMIZATION

- Want to solve $Xw = y$
- If points don't fall perfectly on the line, no solution
- Instead, find the closest approximate solution

$$\min_w \quad || Xw - y ||^2$$

- For some functions, we may not have a closed form solution for the minimum. Find minimum numerically.

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



CONTINUOUS OPTIMIZATION

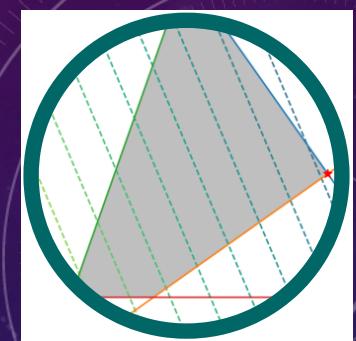
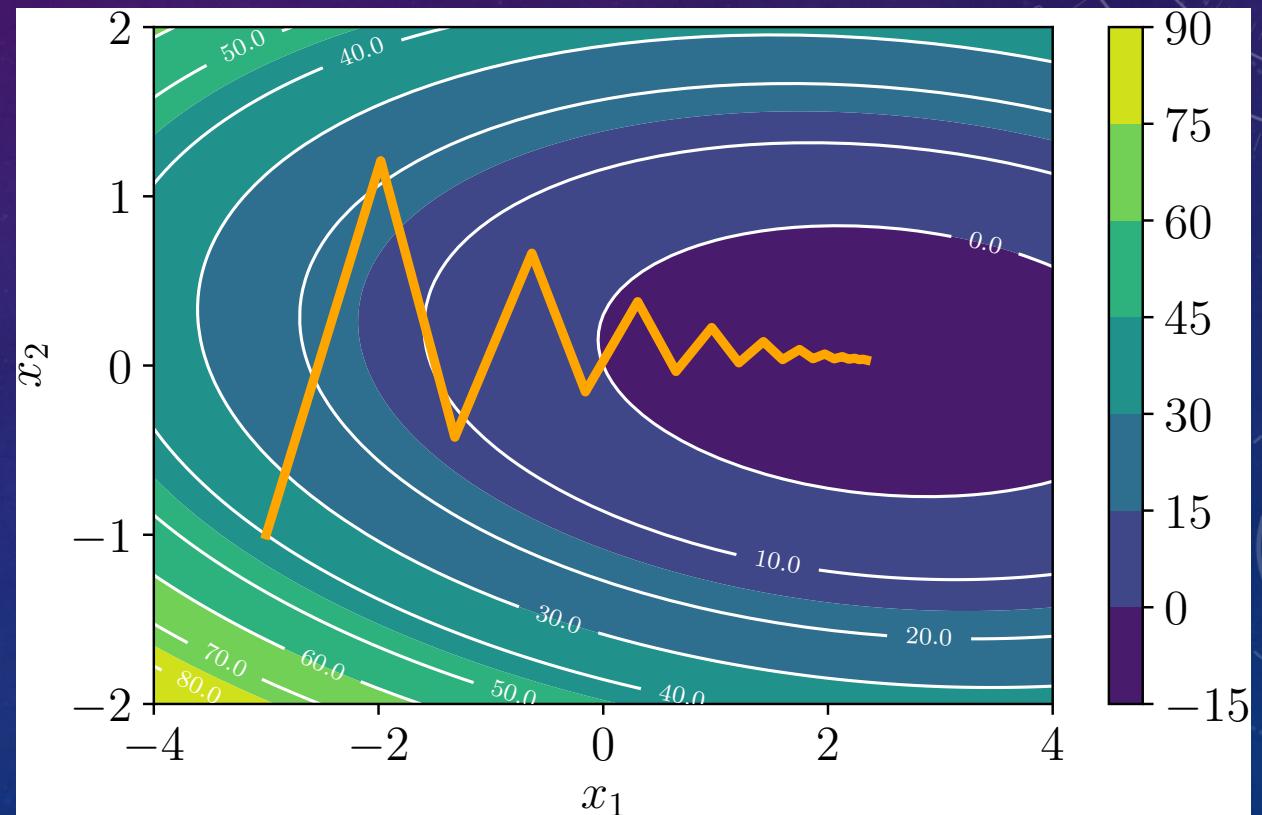
- Objective function

$$\min_{\mathbf{x}} f(\mathbf{x})$$

- Gradient Descent

$$\mathbf{x}_1 = \mathbf{x}_0 - \gamma((\nabla f)(\mathbf{x}_0))^\top$$

- Gradient ∇f
- Step-size γ



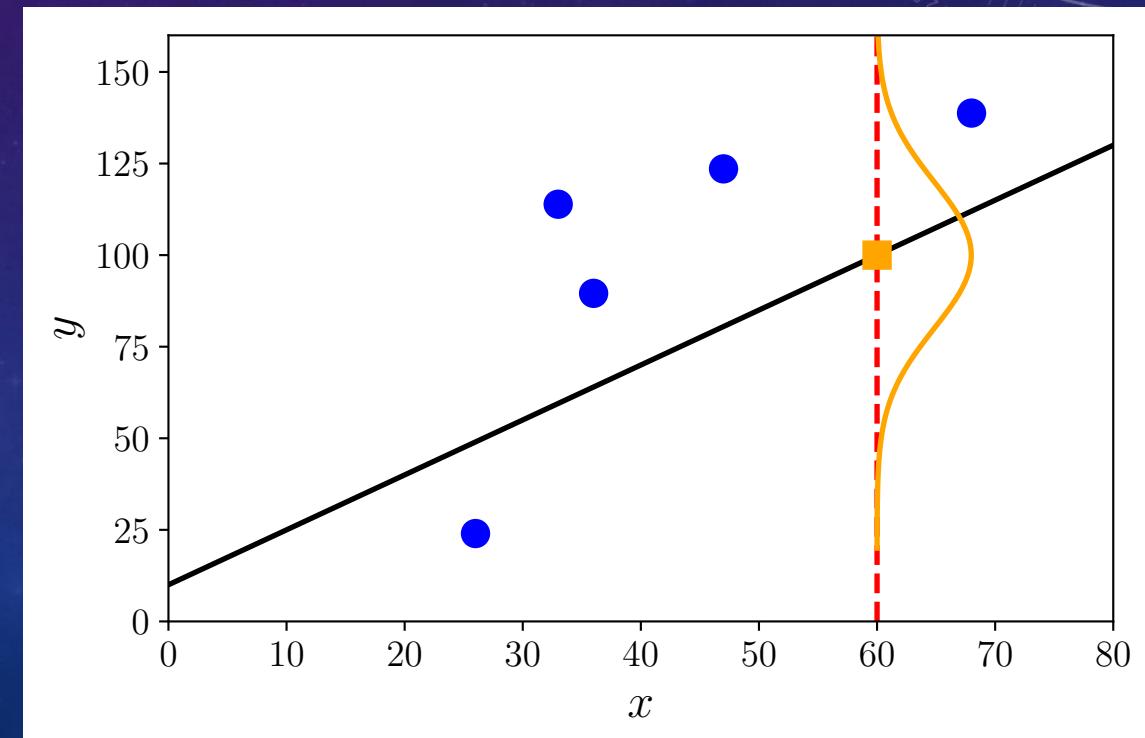
FITTING A LINE – MAXIMUM LIKELIHOOD

- Want to solve $Xw = y$
- If points don't fall perfectly on the line, no solution
- Assume data (X, y) is represented by random variables
- And for a given family of probability densities, compute the maximum likelihood

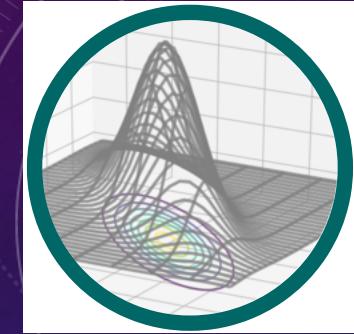
$$\max_w p(y | X, w)$$

- What is the noise model?
- What is the prior?
- What is the predictive uncertainty?

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



PROBABILITY AND DISTRIBUTIONS



- Probability space

- Sample space (Ω), e.g. {hh, ht, th, tt}
- Event space, e.g. one head = {ht, th}
- Probability space, e.g. P(one head) = 0.5

$$S \subseteq \mathcal{T}$$

- Random variable

- Target space, e.g. discrete or real
- Random variable is a function X

$$X : \Omega \rightarrow \mathcal{T}$$

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

Type	“Point probability”	“Interval probability”
Discrete	$P(X = x)$ Probability mass function	Not applicable
Continuous	$p(x)$ Probability density function	$P(X \leq x)$ Cumulative distribution function

PROBABILITY AND DISTRIBUTIONS

Definition 6.1 (Probability Density Function). A function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is called a *probability density function (pdf)* if

1. $\forall \mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) \geq 0$
2. Its integral exists and

$$\int_{\mathbb{R}^D} f(\mathbf{x}) d\mathbf{x} = 1. \quad (6.15)$$

- Distribution (or law) of the random variable

$$P(a \leq X \leq b) = \int_a^b f(x) dx, \quad X : \Omega \rightarrow \mathcal{T}$$

RULES OF PROBABILITY

- Sum rule

$$p(\mathbf{x}) = \begin{cases} \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y}) & \text{if } \mathbf{y} \text{ is discrete} \\ \int_{\mathcal{Y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} & \text{if } \mathbf{y} \text{ is continuous} \end{cases}$$

- Product rule

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})$$

- Bayes' Theorem

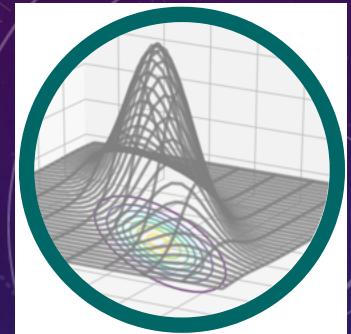
$$p(\mathbf{x} | \mathbf{y}) = \frac{\underbrace{p(\mathbf{y} | \mathbf{x})}_{\text{posterior}} \underbrace{p(\mathbf{x})}_{\text{prior}}}{\underbrace{p(\mathbf{y})}_{\text{evidence}}}$$

Definition 6.3 (Expected Value). The *expected value* of a function $g : \mathbb{R} \rightarrow \mathbb{R}$ of a univariate continuous random variable $X \sim p(x)$ is given by

$$\mathbb{E}_X[g(x)] = \int_{\mathcal{X}} g(x)p(x)dx. \quad (6.28)$$

Definition 6.6 (Covariance (Multivariate)). If we consider two multivariate random variables X and Y with states $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{y} \in \mathbb{R}^E$ respectively, the *covariance* between X and Y is defined as

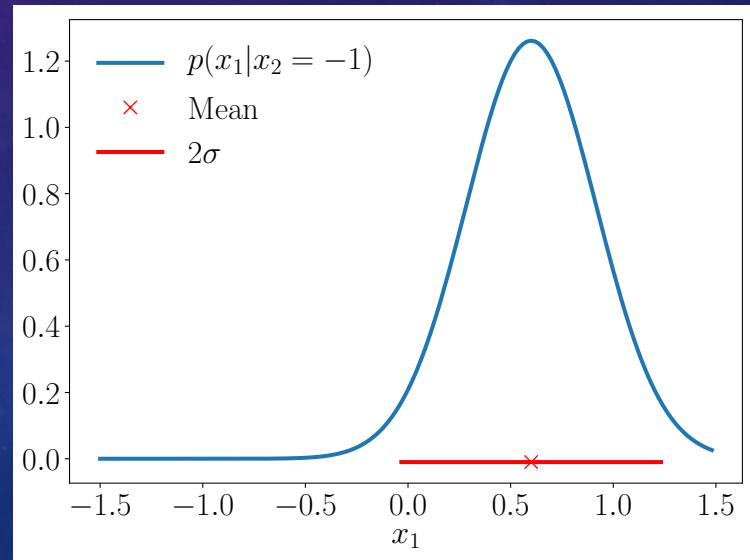
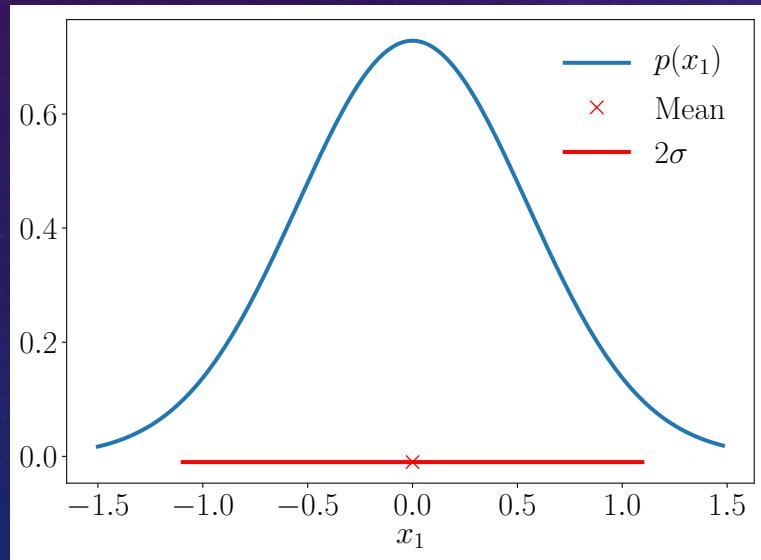
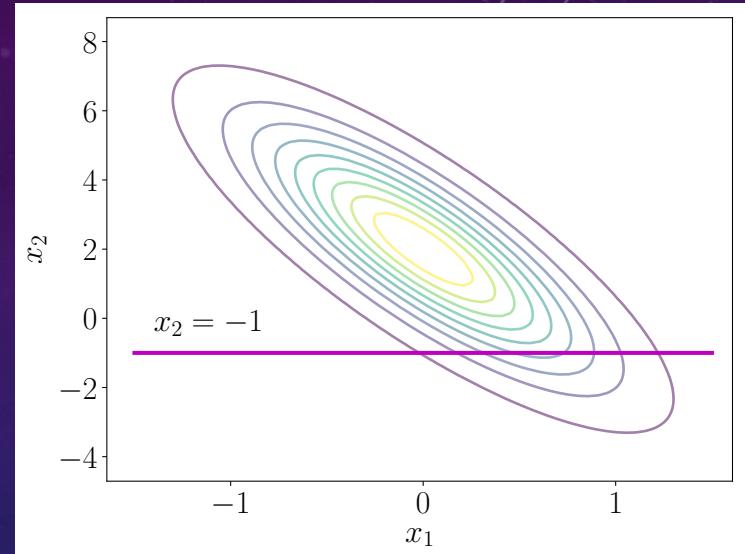
$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]^\top = \text{Cov}[\mathbf{y}, \mathbf{x}]^\top \in \mathbb{R}^{D \times E}. \quad (6.37)$$



GAUSSIAN DISTRIBUTION

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

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$



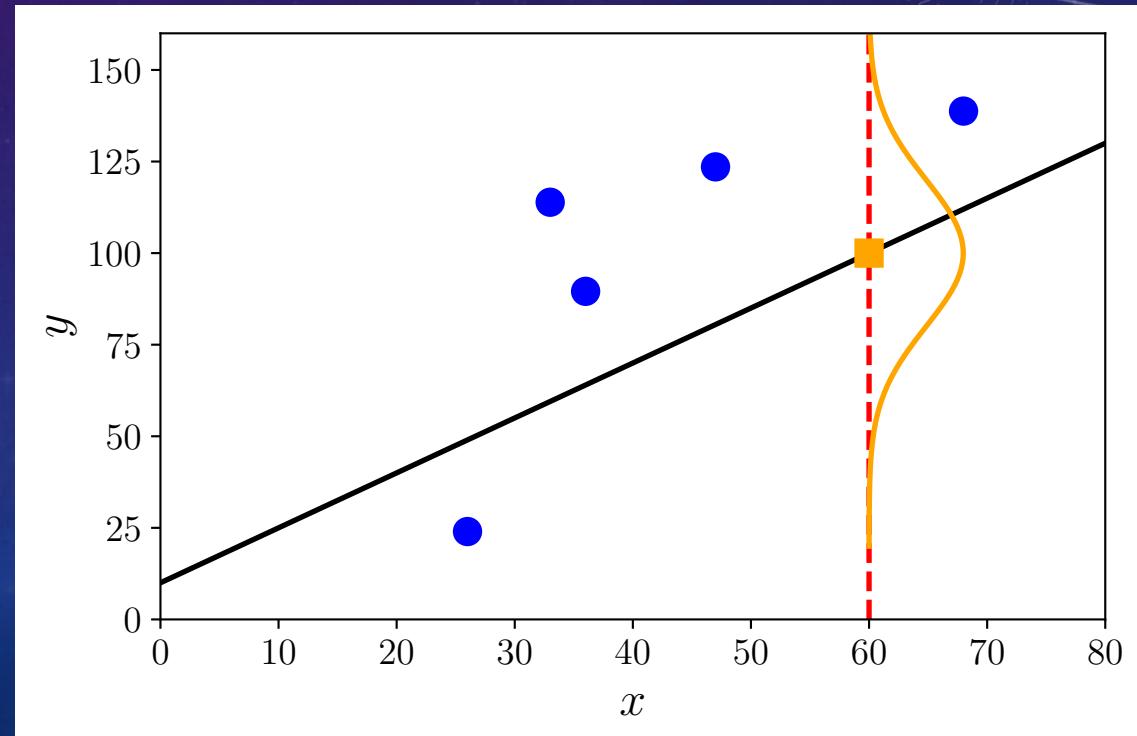
FITTING A LINE – MAXIMUM LIKELIHOOD

- Want to solve $Xw = y$
- If points don't fall perfectly on the line, no solution
- Assume data (X, y) is represented by random variables
- And for a given family of probability densities, compute the maximum likelihood

$$\max_w p(y | X, w)$$

- What is the noise model?
We assume Gaussian noise

Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888

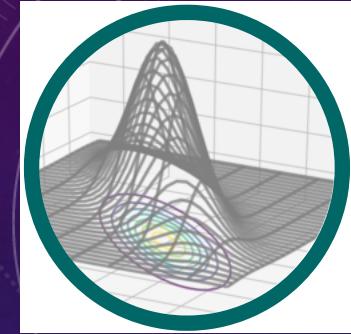


CONJUGACY AND EXPONENTIAL FAMILY

$$p(\mathbf{x} | \mathbf{y}) = \frac{\underbrace{p(\mathbf{y} | \mathbf{x}) p(\mathbf{x})}_{\text{posterior}}}{\underbrace{p(\mathbf{y})}_{\text{evidence}}}$$

Bayes' Theorem

Definition 6.13 (Conjugate Prior). A prior is *conjugate* for the likelihood function if the posterior is of the same form/type as the prior.



An *exponential family* is a family of probability distributions, parameterized by $\boldsymbol{\theta} \in \mathbb{R}^D$, of the form

$$p(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) \exp (\langle \boldsymbol{\theta}, \phi(\mathbf{x}) \rangle - A(\boldsymbol{\theta})) , \quad (6.107)$$

- Natural parameters $\boldsymbol{\theta}$
- Sufficient statistics $\phi(\mathbf{x})$
- Log partition function $A(\boldsymbol{\theta})$

Theorem 6.14 (Fisher-Neyman). [Theorem 6.5 in Lehmann and Casella (1998)] Let X have probability density function $p(x | \boldsymbol{\theta})$. Then the statistics $\phi(\mathbf{x})$ are sufficient for $\boldsymbol{\theta}$ if and only if $p(x | \boldsymbol{\theta})$ can be written in the form

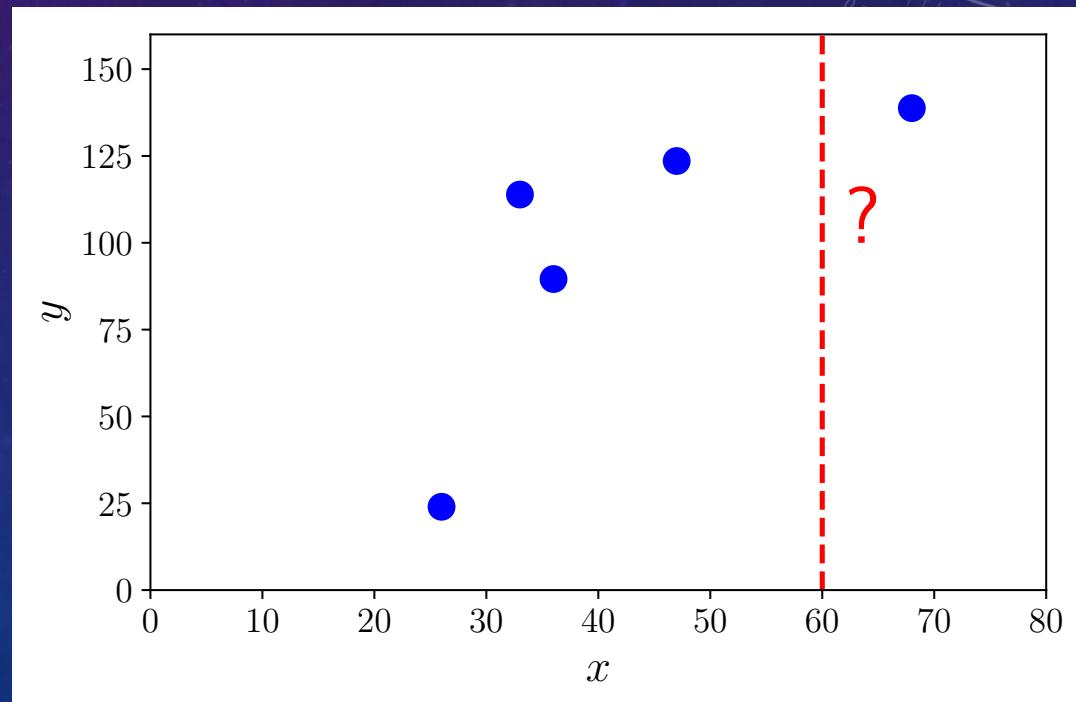
$$p(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) g_{\boldsymbol{\theta}}(\phi(\mathbf{x})) , \quad (6.106)$$

where $h(\mathbf{x})$ is a distribution independent of $\boldsymbol{\theta}$ and $g_{\boldsymbol{\theta}}$ captures all the dependence on $\boldsymbol{\theta}$ via sufficient statistics $\phi(\mathbf{x})$.

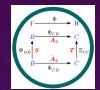
MACHINE LEARNING IS ABOUT PREDICTION

- Predict salary (y) from age (x)
- The values of y for the training data is not the main focus
- We are interested in generalization error:
 - What is the error we make on unseen data?
- Do not train on the test set

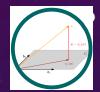
Age	Annual Salary (in thousands)
36	89.563
47	123.543
26	23.989
68	138.769
33	113.888



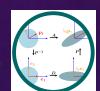
SUMMARY: MATH FOR ML



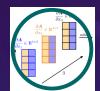
Chap. 2: Linear Algebra
vector space, linear maps, affine space



Chap. 3: Analytic Geometry
inner products, distances, orthogonality



Chap. 4: Matrix Decompositions
Cholesky, eigenvalue decomposition, singular value decomposition



Chap. 5: Vector Calculus
gradient is a row vector, chain rule



Chap. 6: Probability and Distributions
random variable, distribution, Bayes rule, expectation



Chap. 7: Continuous Optimization
gradient descent, convex duality



MATHEMATICS FOR MACHINE LEARNING



... and 5 chapters of ML

mml-book.com