# Lecture 19
# Elements of Machine Learning

Basic terminology, concepts, and methods; linear models

*Prof. David A. Kofke*
*CE 500 – Modeling Potential-Energy Surfaces*
*Department of Chemical & Biological Engineering*
*University at Buffalo*

**University at Buffalo**
The State University of New York

- New skill                          d through
  observatio

- Most algo                          s

- Core elen
  – Statistic
  – Artifici

- What nev

  – The inte                    s of data

  – Motivates                                es of algorithms                    r useful purposes

TO COMPLETE YOUR REGISTRATION, PLEASE TELL US WHETHER OR NOT THIS IMAGE CONTAINS A STOP SIGN:

NO    YES

ANSWER QUICKLY—OUR SELF-DRIVING CAR IS ALMOST AT THE INTERSECTION.

SO MUCH OF "AI" IS JUST FIGURING OUT WAYS TO OFFLOAD WORK ONTO RANDOM STRANGERS.
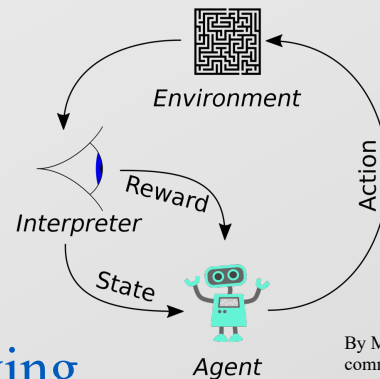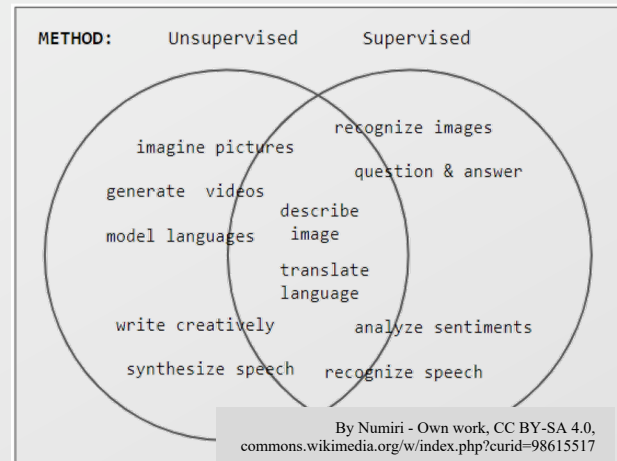
Randall Munroe
xkcd.com/1897/

# *Machine learning* (ML) is the study of algorithms that auto-improve via experience and data

- New skills and/or better performance is acquired through observation and trial & error

- Most algorithms are based in statistical concepts

- Core elements have been around for decades
  - Statistical inference methods
  - Artificial intelligence

- What new?  Data!
  - The internet and social media produces huge amounts of data
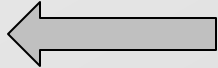  - Motivates development of algorithms to harness for useful purposes

- Supervised
  - {input, output} data pairs are provided, and goal is to provide correct output values for new input data

- Unsupervised
  - No output labels are provided with data; rather algorithm seeks to find patterns in it

- Reinforcement
  - Agent explores actions guided by rewards
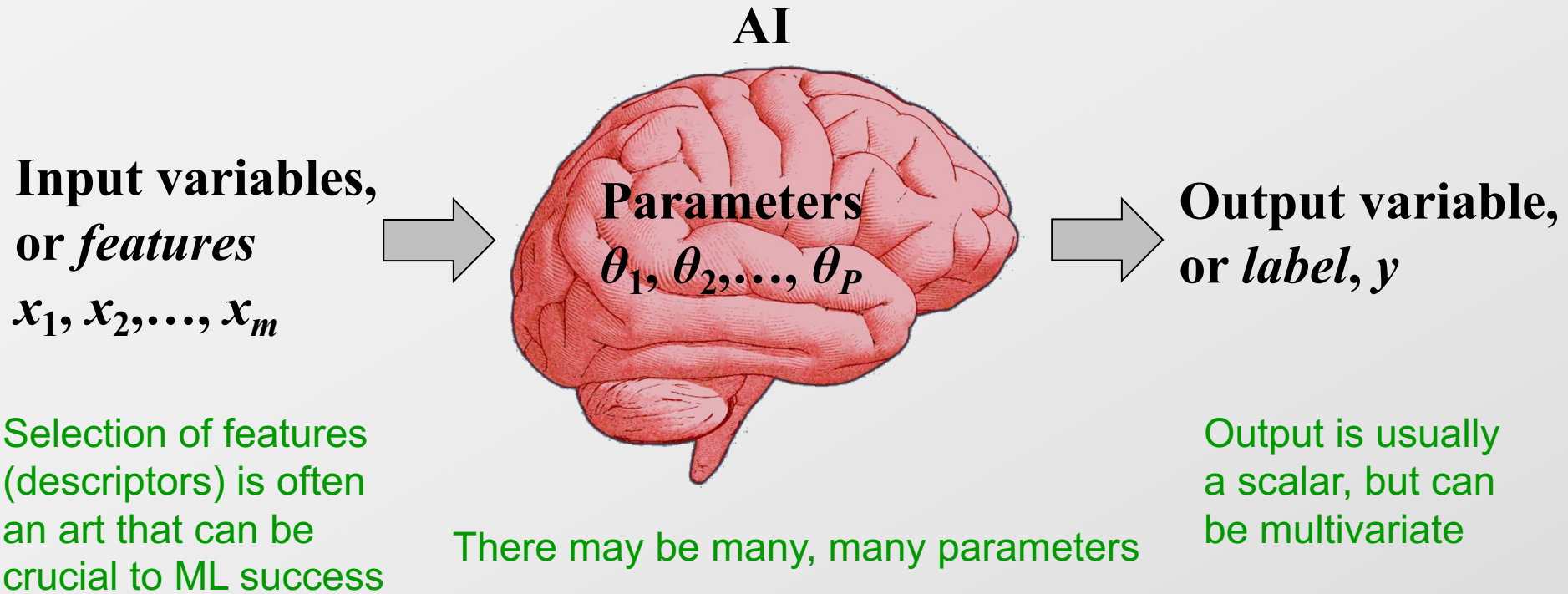  - E.g., games, robot control, autonomous driving



METHOD:  Unsupervised    Supervised

imagine pictures    recognize images
generate videos    question & answer
model languages    describe image
    translate language
write creatively    analyze sentiments
synthesize speech    recognize speech

By Numiri - Own work, CC BY-SA 4.0, commons.wikimedia.org/w/index.php?curid=98615517



Environment

Action

Reward

Interpreter

State

Agent

By Megajuice - Own work, CC0, commons.wikimedia.org/w/index.php?curid=57895741

4

- Classification aims to identify the discrete category for new input data

  - E.g., whether a new molecule is toxic vs. non-toxic

- Regression aims to estimate the value of a continuous variable given new input data

  - E.g., $pK_a$, atomization energy, redox potential for a new molecule
  - Potential energy for a configuration of molecules

# The operation of a supervised ML algorithm is governed by a set of numeric parameters

**AI**

**Input variables, or *features***
$x_1, x_2, \ldots, x_m$

**Parameters**
$\theta_1, \theta_2, \ldots, \theta_P$

**Output variable, or *label*, $y$**

Selection of features (descriptors) is often an art that can be crucial to ML success

There may be many, many parameters

Output is usually a scalar, but can be multivariate

# *Training* is the process of establishing parameter values, by minimizing a *cost function*

- The training set is a collection of $(\mathbf{x}, y)$ pairs

- A cost function (aka *loss*, *error*) characterizes the error in the ML estimate of the output values relative to the given ones

  – L1 norm: $\displaystyle\sum_{i=1}^{n} \left| \hat{y}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) - y^{(i)} \right| = ||\hat{\mathbf{y}} - \mathbf{y}||_1$

  – L2 norm: $\displaystyle\sum_{i=1}^{n} \left( \hat{y}(\mathbf{x}^{(i)}; \boldsymbol{\theta}) - y^{(i)} \right)^2 \equiv ||\hat{\mathbf{y}} - \mathbf{y}||_2^2$

- The training process attempts to minimize the cost function through manipulation of the parameters $\boldsymbol{\theta}$

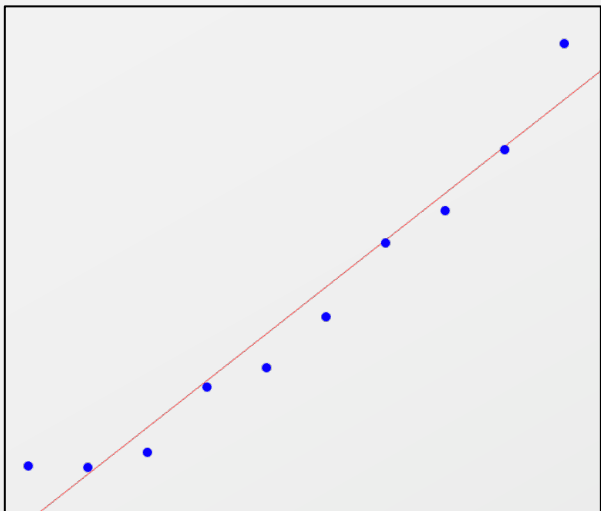# Some data processing may be performed before starting the training and application

- Cleaning
  - *e.g.*, filling in missing values

- Standardization or other transformation
  - *e.g.,* scaling to zero mean and unit standard deviation
  - Often this makes fitting more generic and easier, without irreversibly changing the data

# Available data should be split into training, validation, and test sets

| Training set | Validation set | Test set |
|---|---|---|

- *Training set* to determine the ML model parameters

- *Validation set* to adjust hyperparameters and avoid overfitting
  - Hyperparameters define structure of ML model or guide training
  - Validation may be added to training once hyperparameters are set

- *Test set* to assess the ML model

  - It should play no part in training or validation

- Data should be distributed at random among the sets
  - 65:15:20 distribution of training:validation:test is typical

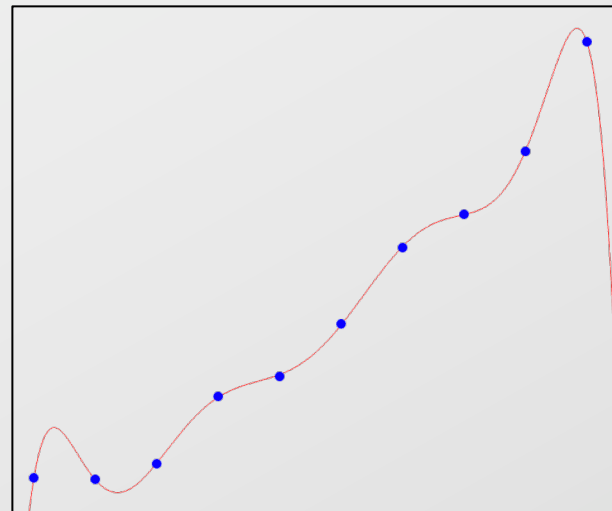# Both over- and underfitting are bad



### Underfitting
Meaningful relationships in datasets are not learned

training error: high
validation error: high

### Appropriate fitting

training error: low
validation error: low

### Overfitting
Noise in training data is learned, and does not generalize

training error: low
validation error: high

Look at the values of the parameters for each of these fits $\sum a_i x^i$

$a_0 = -5.7$
$a_1 = 11.3$

$a_0 = 1.4$
$a_1 = 3.0$
$a_2 = 1.5$

$a_0 = -12,010$
$a_1 = 110,100$
$a_2 = -432,100$
$a_3 = 947,700$
$a_4 = -1,272,000$
*etc.* (up to $a_8$)

→ Overfitting often achieved using very large parameter values

Underfitting
Meaningful relationships in datasets are not learned

Appropriate fitting

Overfitting
Noise in training data is learned, and does not generalize

training error: high
validation error: high

training error: low
validation error: low

training error: low
validation error: high

# *Regularization* calibrates ML models to prevent underfitting or overfitting

- One approach: add a cost-function penalty for large parameters
  - L1 regularization or *lasso* regression:
    $$\text{cost} = \text{error} + \lambda \sum_i |\theta_i|$$

    - Favors sparsity of coefficients, making some exactly 0
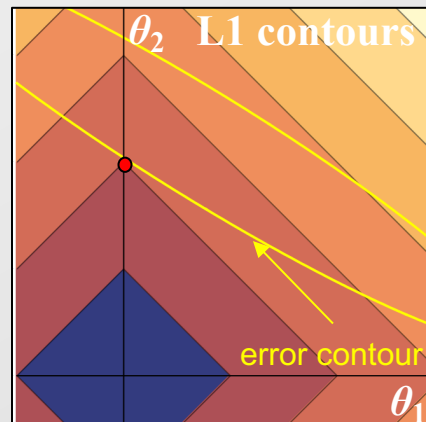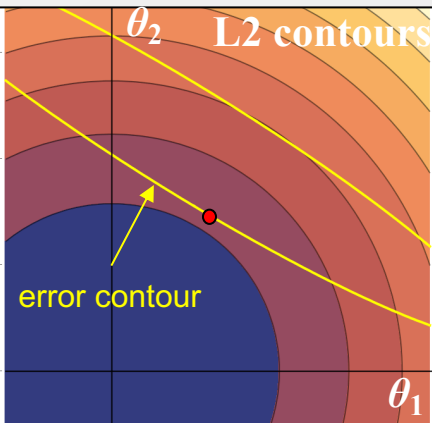    - A tool for feature selection

    $\theta_2$ **L1 contours**

    error contour

    $\theta_1$

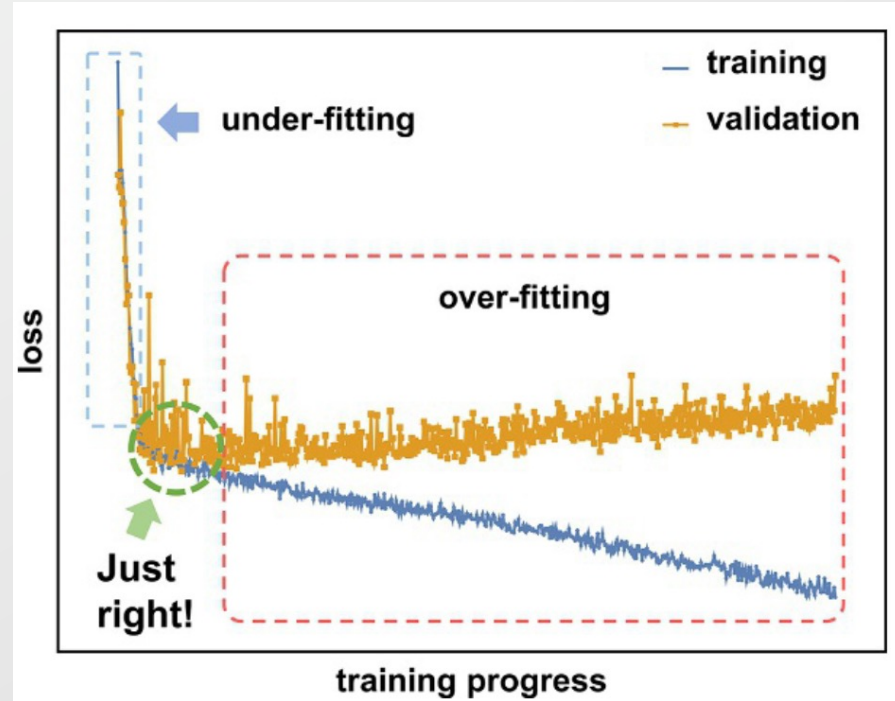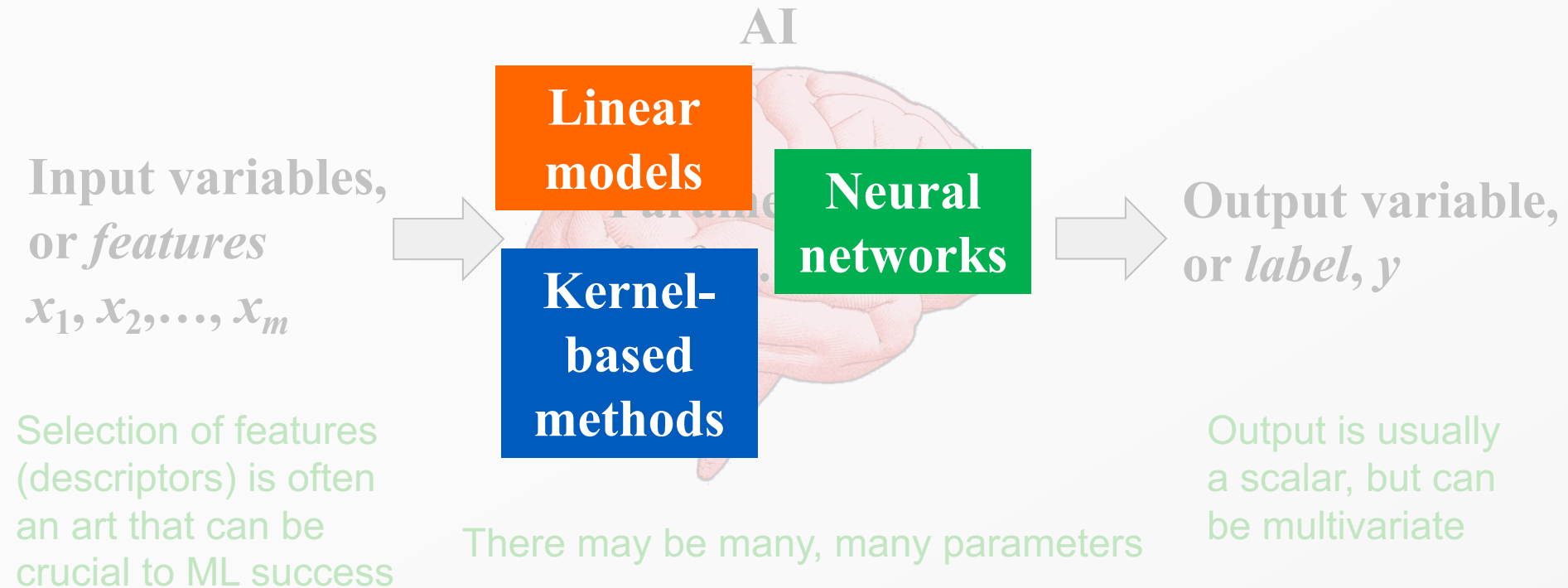    λ is an example of a hyperparameter

    For L1 regularization, intersection with minimum contour is likely to happen at a vertex of penalty function

  - L2 regularization, or *ridge* regression:
    $$\text{cost} = \text{error} + \lambda \sum_i \theta_i^2$$

    $\theta_2$ **L2 contours**

    error contour

    $\theta_1$

    Overall compression toward smaller coefficients
    Also known as *Tikhonov regularization*

# *Early stopping* is another regularization method

- Perform parameter optimization on training set

- Occasionally evaluate error using validation set

- Where validation error begins to increase, halt optimization



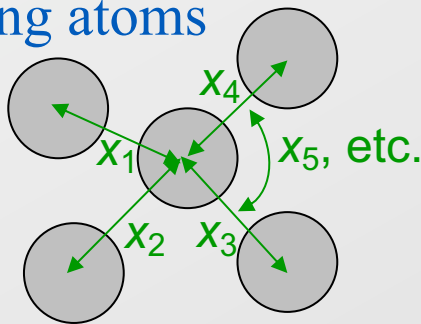Dral et al., doi.org/10.1016/B978-0-323-90049-2.00011-1

# A model is *linear* if it has linear dependence on its parameters (not that it fits using linear functions)

- Energy given as a sum of one-body energies

$$E = \sum_i^{N_{\text{atoms}}} E_i$$

– Atom energies are, in turn, given via a set of descriptors that depend on positions of neighboring atoms



$x_4$

$x_1$ $x_5$, etc.

$x_2$ $x_3$

- General linear model has a simple form: $\mathbf{y} = \mathbf{X}\widehat{\mathbf{w}}$

$$\begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{pmatrix} = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_p^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_p^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \cdots & x_p^{(n)} \end{pmatrix} \begin{pmatrix} \widehat{w}_1 \\ \widehat{w}_2 \\ \vdots \\ \widehat{w}_p \end{pmatrix}$$

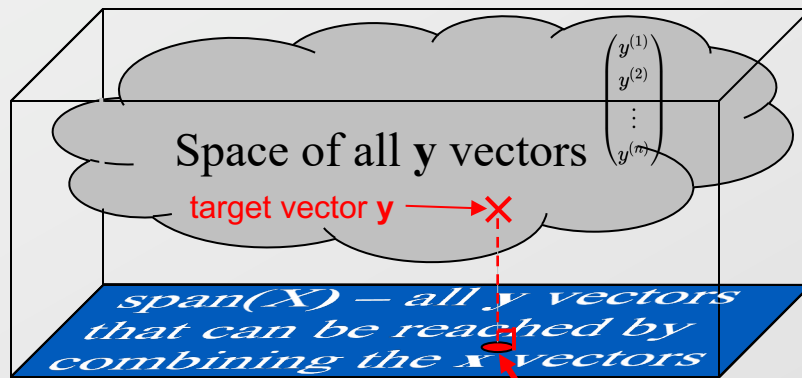*n* {**x**,*y*} observation pairs

*p* features and parameters

16

$$\widehat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w} \in \mathcal{R}^p} \sum_{i=1}^n \left( y^{(i)} - \sum_j^p x_j^{(i)} w_j \right)^2 = \operatorname*{argmin}_{\mathbf{w} \in \mathcal{R}^p} ||\mathbf{y} - \mathbf{X}\mathbf{w}||_2^2$$

- Given a set of $x_j$ vectors, how can we combine them to get as close to the **y** vector as possible? Think geometrically.

$$\begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_p^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_p^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \cdots & x_p^{(n)} \end{pmatrix}$$

$\uparrow$ $\uparrow$ $\uparrow$

$\mathbf{x}_1$ $\mathbf{x}_2$ $\ldots$ $\mathbf{x}_p$

$\mathbf{x}_j$ vectors



Space of all **y** vectors

$\begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{pmatrix}$

target vector **y**

*span(X) — all y vectors that can be reached by combining the x vectors*

closest vector in span(**X**): $\widehat{\mathbf{y}} = \mathbf{X}\widehat{\mathbf{w}}$

$\mathbf{y} - \widehat{\mathbf{y}}$ is orthogonal to all $\mathbf{x}_j$:

$\mathbf{X}^{\mathrm{T}}(\mathbf{y} - \widehat{\mathbf{y}}) = 0$

$\mathbf{X}^{\mathrm{T}}\mathbf{X}\widehat{\mathbf{w}} = \mathbf{X}^{\mathrm{T}}\mathbf{y}$

$\boxed{\widehat{\mathbf{w}} = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}}$

Direct solution for optimum is obtained (if enough data)

# Regularization is needed if the parameters exceed the number of data (XᵀX not invertible)

- Where $\mathbf{w}$ doesn't have a unique solution, its evaluation is arbitrary to a degree, and prediction performance will suffer

  – Situation is likely where $p > n$

  – Regularization can alleviate this

- Ridge regression

$$\widehat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w} \in \mathcal{R}^p} \left( \tfrac{1}{2} ||\mathbf{y} - \mathbf{X}\mathbf{w}||_2^2 + \lambda ||\mathbf{w}||_2^2 \right)$$

- Still a quadratic form; analytic minimum

$$\widehat{\mathbf{w}} = (\mathbf{X}^{\mathrm{T}}\mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^{\mathrm{T}}\mathbf{y}$$

  – Inverse will exist for nonzero $\lambda$

$$\mathbf{X} = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_p^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \cdots & x_p^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \cdots & x_p^{(n)} \end{pmatrix}$$

$$\widehat{\mathbf{w}} = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

$$\widehat{\mathbf{w}} = \underset{\mathbf{w} \in \mathcal{R}^p}{\operatorname{argmin}} \left( \tfrac{1}{2} ||\mathbf{y} - \mathbf{X}\mathbf{w}||_2^2 + \lambda ||\mathbf{w}||_{\boxed{1}} \right)$$
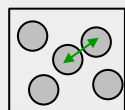
- Not a quadratic form, so more complicated to minimize

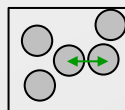- Consider a single-feature example  $\widehat{w} = \underset{w \in \mathcal{R}}{\operatorname{argmin}} f(w)$

$$f(w) = \tfrac{1}{2}\left( w^2 ||\mathbf{x}||_2^2 + ||\mathbf{y}||_2^2 \right) - w\mathbf{x}^{\mathrm{T}}\mathbf{y} + \lambda |w|$$

Example

$x_1 \equiv r_{ij}^{\min}$

$x_1^{(1)} = 0.5$
$y^{(1)} = -2$

$x_1^{(2)} = 0.2$
$y^{(2)} = -7$

$x_1^{(3)} = 0.7$
$y^{(3)} = -3$

Piecewise quadratic

q1   q2

f(w)

λ=10

λ=1

In general, $\widehat{w} = 0$
for $|\mathbf{x}_1^{\mathrm{T}}\mathbf{y}| \le \lambda$

$\mathbf{x}_1$      $\mathbf{y}$

$\begin{pmatrix} 0.5 \\ 0.2 \\ 0.7 \end{pmatrix}$  $\begin{pmatrix} -2 \\ -7 \\ -3 \end{pmatrix}$

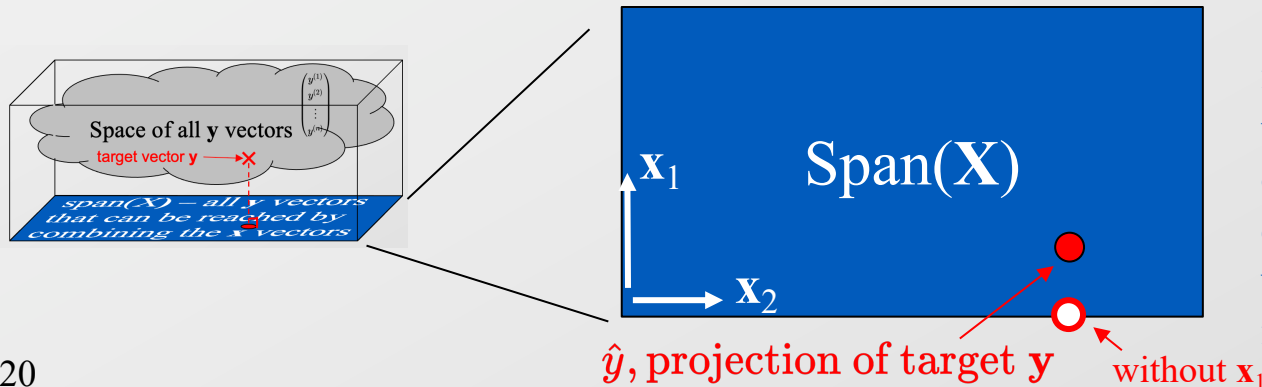$\mathbf{x}_1^{\mathrm{T}}\mathbf{y} = -4.5$

$\widehat{w}$

19

# General lasso case is also piecewise quadratic, but with more ($2^p$) pieces

$$f(w) = \tfrac{1}{2}\mathbf{w}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}\mathbf{X}\mathbf{w} - \mathbf{w}^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}\mathbf{y} + \tfrac{1}{2}\|\mathbf{y}\|_2^2 + \lambda\|\mathbf{w}\|_1$$

- Solution indicates
  $$\left|\mathbf{x}_k^{\mathrm{T}}(\mathbf{y} - \mathbf{X}\widehat{\mathbf{w}})\right| \le \lambda \qquad w_k = 0$$
  $$\left|\mathbf{x}_k^{\mathrm{T}}(\mathbf{y} - \underbrace{\mathbf{X}\widehat{\mathbf{w}}}_{\hat{y}})\right| > \lambda \qquad \mathbf{x}_k^{\mathrm{T}}(\mathbf{y} - \mathbf{X}\widehat{\mathbf{w}}) = \lambda\,\mathrm{sgn}(w_k)$$
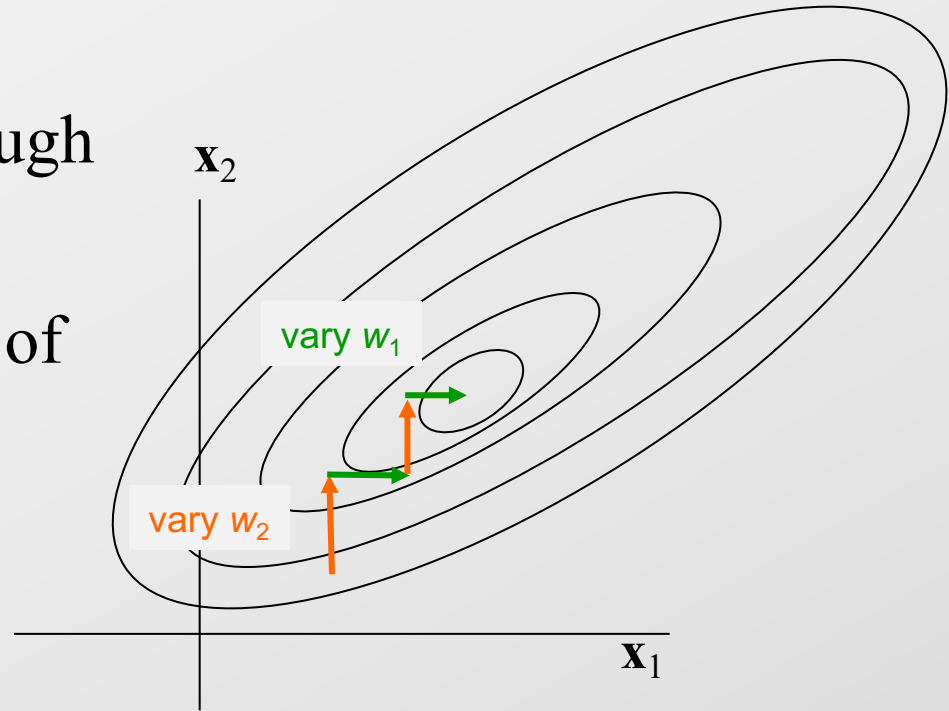
  – Not a closed-form solution, but it gives some geometric insight



if $\mathbf{x}_1^{\mathrm{T}}(\mathbf{y} - \widehat{\mathbf{y}}) \le \lambda$ then feature #1 is almost orthogonal to difference and doesn't add enough to be worthwhile to keep, so its weight is zeroed out
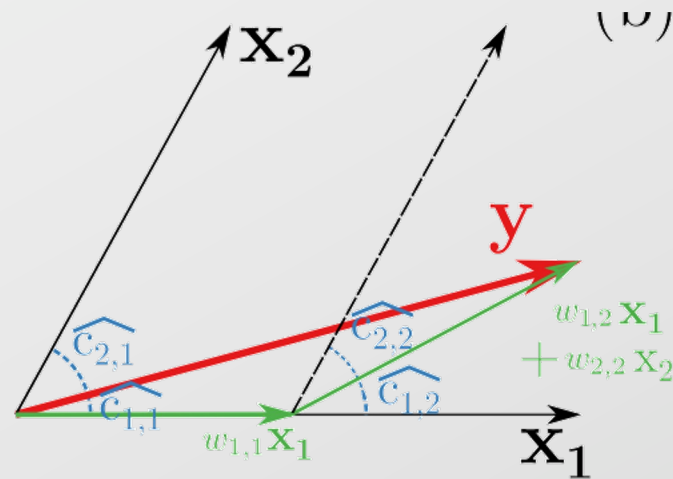
$\hat{y}$, projection of target $\mathbf{y}$

without $\mathbf{x}_1$

# Coordinate descent is a popular algorithm for finding lasso optimum

- Pick initial guess $\mathbf{w}$

- In iteration $i$, proceed through all $w_k$ in random order

- For each $k$, find minimum of $f(\mathbf{w})$ by varying only $w_k$

- Repeat for next $i$ to convergence

- Start with all $w_k = 0$

- Select $k$ for which $\mathbf{x}_k^T\mathbf{y}$ is largest

- Increase $w_k$ until another descriptor $j$ has larger correlation with residual: $\mathbf{x}_j^T(\mathbf{y}-w_k\mathbf{x}_k)$

- Optimize $w_k$ and $w_j$ in direction equiangular between $\mathbf{x}_k$ and $\mathbf{x}_j$ until a 3$^{\text{rd}}$ is more correlated with new residual

- 



Tallec et al.

- Eugen Hruska and Fang Liu, Chapter 6, Machine learning: An overview. In *Quantum Chemistry in the Age of Machine Learning*.
  - Posted on UBLearns

- Gauthier Tallec, Gaétan Laurens, Owen Fresse-Colson, and Julien Lam, Chapter 11, Potentials based on linear models. In *Quantum Chemistry in the Age of Machine Learning*.
  - Posted on UBLearns