

## Keywords

- what is ML
- test/train/validation
- linear algebra
- regularization

elastic net

$$\left\{ \begin{array}{l} q=1 \text{ LASSO } \sum e^2 + \alpha \frac{1}{2} \sum |\theta_i| \\ q=2 \text{ Ridge / Tikhonov } \dots + \alpha \frac{1}{2} \sum |\theta_i|^2 \\ q=0 \text{ minmax} \end{array} \right.$$

$x \in \mathbb{R}$

$$\min \|x - b\|_q$$

$\left\{ \begin{array}{l} \text{median} \\ \text{mean} \\ \text{midrange} \end{array} \right.$

- missing data

- examples

- norms

- maximum likelihood

- feature engineering

• bayesian linear regression  
• logistic regression  
(cross-entropy)

$$y = f(w^T \phi)$$

↑ activation function  
 $f^{-1}$  = link function

## Part 1

- what is ML (tabular)
- examples, relational data
- linear regression
- norms
- feature engineering
- test/train/validation
- numpy, pandas, mpl, sklearn

## Part 2

- Recap
- Regularization
- categorical data, missing data
- generalized linear models  
(poisson, logistic regression)
- sklearn pipelines
- tree models

## Sources:

kap 6 "Approximation and fitting" Boyd

kap 1, 2, 3, 4, 6, 7 "Hands on ML" Géron

kap 3, 4 "Pattern Recognition" Bishop

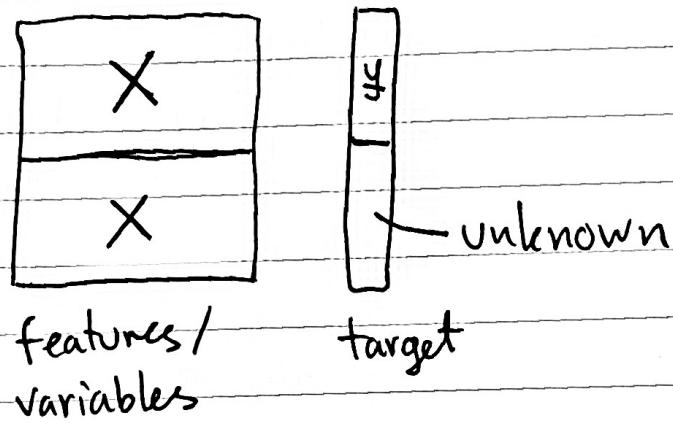
kap 3, 4, 9, 10 "Elements" Hastie

# Fagdag Sonat

What is machine learning?

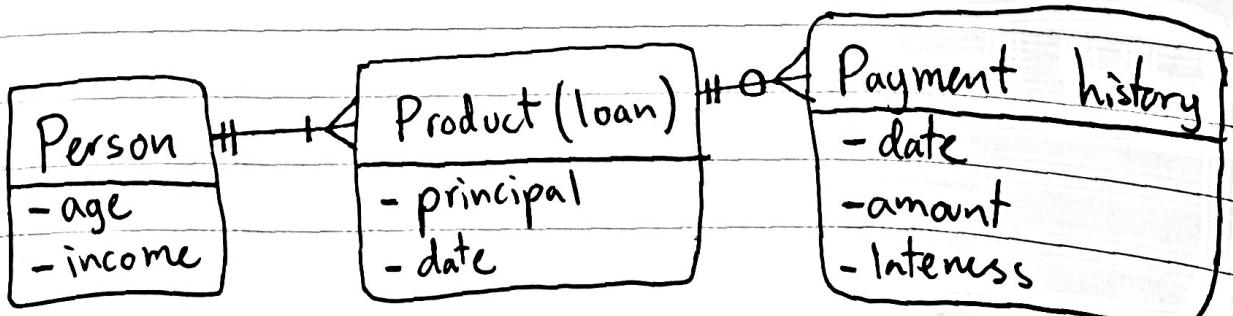
$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

- Pairs  $(x_n, y_n)$  in a table



## Examples

- Doctor knows  $x_n = (\text{age}, \text{weight}, \text{height})$  and wants to learn  $y_n = \text{blood pressure}$
- Houses with  $x_n = (\text{sqm}, \text{age}, \text{floor})$ , investor wants to know  $y_n = \text{sales price}$
- Probability of default given products and payments (relational, tree-like structure)



## Is it really machine learning?

- statistical inference ("which variables affect blood pressure")
- optimization ("find best solution / configuration")  
"which house is undervalued"
- agent navigating search space ("play chess")

ML = function approximation for the purpose of out-of-sample prediction

## Linear regression

$$\hat{y}_n = \underline{w}^T \underline{x}_n = \sum_i w_i x_{ni} \quad \hat{y} = \underline{X} \underline{w}$$

$$\widehat{\text{blood-pressure}} = w_1 \cdot \text{age} + w_2 \cdot \text{weight} + w_3 \cdot \text{height}$$

$$\underline{e} = (\underline{y} - \hat{\underline{y}}) = (\underline{y} - \underline{X} \underline{w})$$

$\underline{e}$  = vector with n error terms

OLS : ordinary least squares  
minimize  $\underline{e}^T \underline{e}$

$$\|\underline{e}\|_2^2 = \underline{e}^T \underline{e} = (\underline{y} - \underline{X} \underline{w})^T (\underline{y} - \underline{X} \underline{w}) = \sum_n e_n^2 = \sum_n [y_n - \underline{x}_n^T \underline{w}]^2$$

Convex optimization problem, easily solved.

→ All errors are not the same - norms

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{1/p}$$

p-norms, or  $\ell$ -norms

$p=1$  Manhattan distance

$p=2$  Euclidean distance

$p=\infty$  Max-norm

Consider

$$\min_x \|x - a\|$$

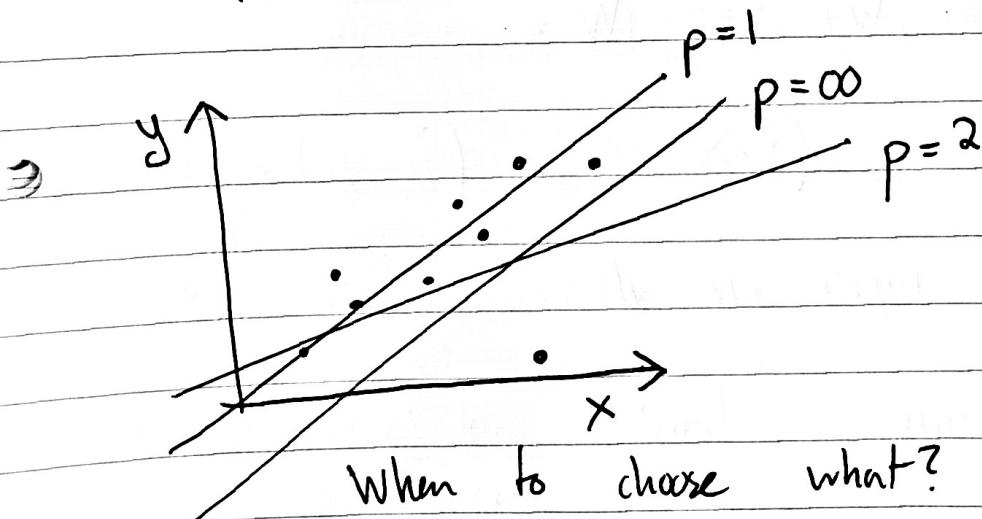
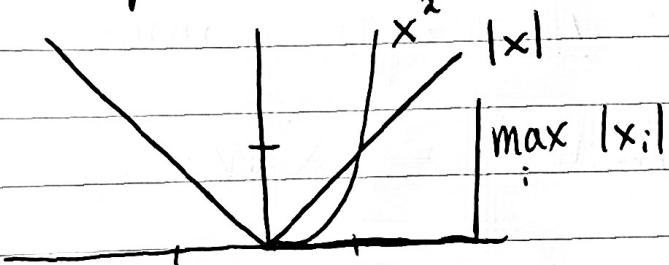
$x$

median, mean, mid-range

→  $p=2$   $e=(2,2)$  better than  $(4,0)$ , "fair"

$p=1$  less sensitive to outliers

$p=\infty$  only cares about worst case (min max)



When to choose what?

$p=2$  distribute errors, few outliers

$p=1$  errors tied to economics, serious outliers

$p=\infty$  make worst case as ok as possible

## → Feature engineering

$\text{blood-pressure} = w_1 \cdot \text{age} + w_2 \cdot \text{weight} + w_3 \cdot \text{height} + w_4 \cdot \frac{\text{height}}{\text{height}}$

features in data      new, engineered feature

$\phi$  is the feature-mapping

$\hat{y}_n = \underline{w^T} \phi(x_n)$  instead of  $\hat{y}_n = \underline{w^T} x_n$

Example 1 Altman's financial ratios for predicting bankruptcy in 1968.

Example 2 Using "income" is a bad idea. Feature will "drift" because of inflation.

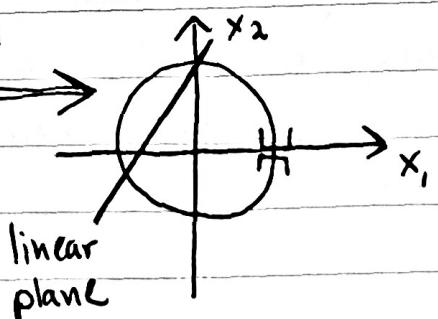
Example 3 If a variable is log-normal, taking logs often helps.

Example 4 Predicting ice-cream sales using time as a variable.

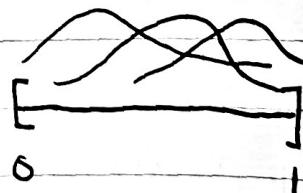
$t \in [0, 1]$  want

$$x_1 = \sin(2\pi t)$$

$$x_2 = \cos(2\pi t)$$



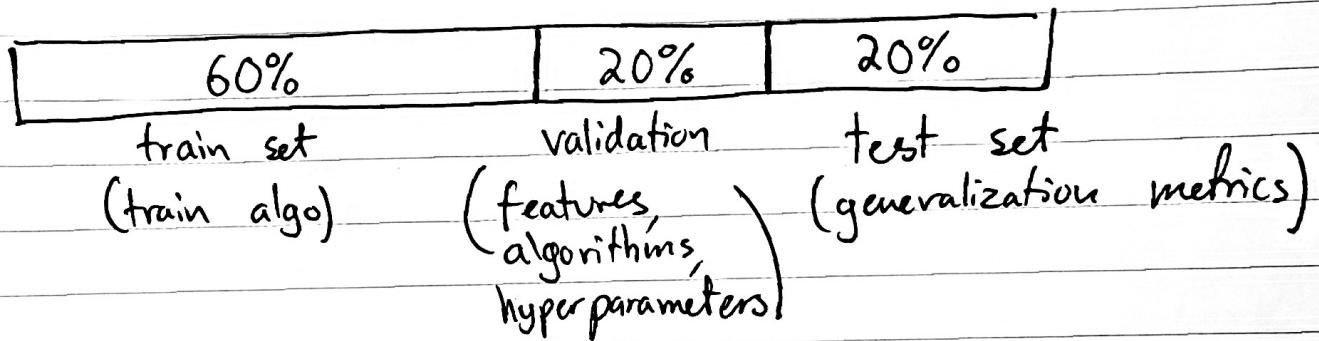
$$y = w_1 \sin(2\pi t) + w_2 \cos(2\pi t)$$



## ⇒ Feature engineering tips

- ask domain experts
- prefer dimensionless variables
- try logs, exp, polynomial combinations
- use relational structures, aggregate:
  - mean, variance, length, max-min, etc...

## ⇒ Train, validate, test sets



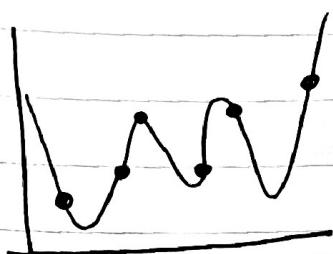
- or use k-folds cross validation

## ⇒ Python and ML

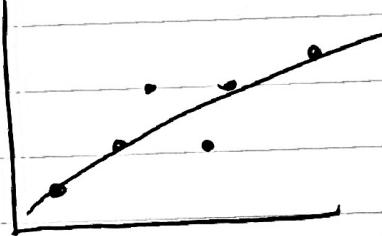
- pandas
- scikit-learn

## Regularization

- penalize complex models to avoid overfitting (ex: little data,



overfitting



regularized

$$\underset{\underline{w}}{\text{minimize}} \quad \left( \underbrace{\|\underline{y} - \underline{X}\underline{w}\|_2^2}_{\text{model fit error}}, \underbrace{\|\underline{w}\|_p^2}_{\text{model complexity}} \right)$$

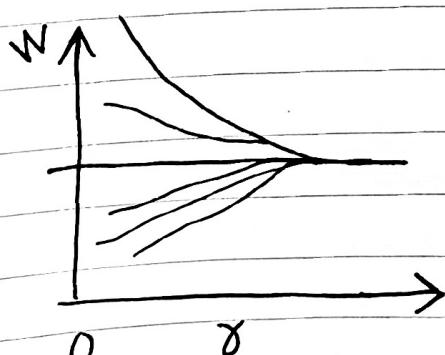
multi-objective problem,  
scalarized as

Ridge  $p=2$

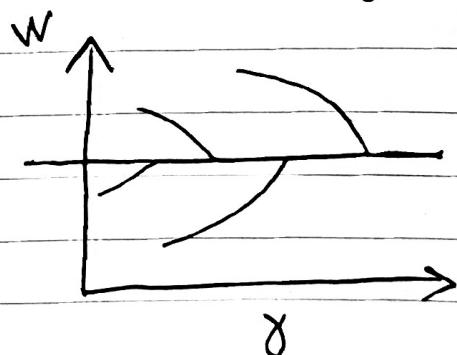
$$\min_{\underline{w}} \|\underline{y} - \underline{X}\underline{w}\|_2^2 + \gamma \|\underline{w}\|_2^2$$

Lasso  $p=1$

$$\min_{\underline{w}} \|\underline{y} - \underline{X}\underline{w}\|_2^2 + \gamma \|\underline{w}\|_1$$



Regularization strength



- Data types
  - categorical ("color") [nominal] (no order)
  - integer ("age") [ordinal] (rank order)
  - float/real ("pressure") [interval scale]  
[ratio scale]

### Examples

- grades {"A", "B", "C", ...}
- post codes {"4016", "4017", ...}

### • One hot encoding

color	color_blue	color_red	color_yellow
blue	1	0	0
red	0	1	0
red	0	1	0
yellow	0	0	1

⇒

- can use most common entries if there are many.

### • Missing data

- what does missing mean? unmeasured, not existing
- is data missing at random, or does it mean anything for y? ("age" vs. "car" missing)

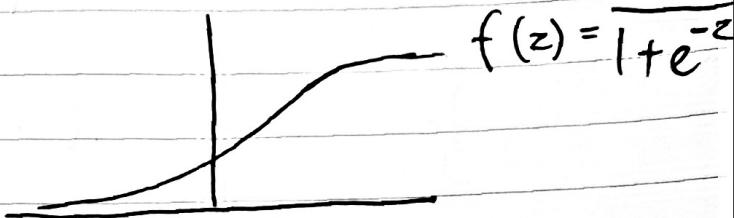
### Strategies

- set to zero ("car age")
- set to median of column/variable
- drop rows with missing data
- KNN / multiple imputations

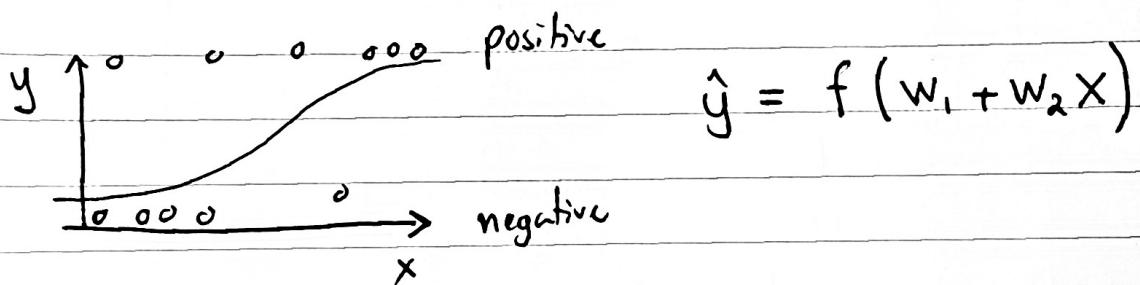
## Generalized linear models

$$\hat{y}_n = f(\underline{w}^\top \varphi(\underline{x}_n))$$

$$f(z) = \frac{1}{1 + e^{-z}}$$



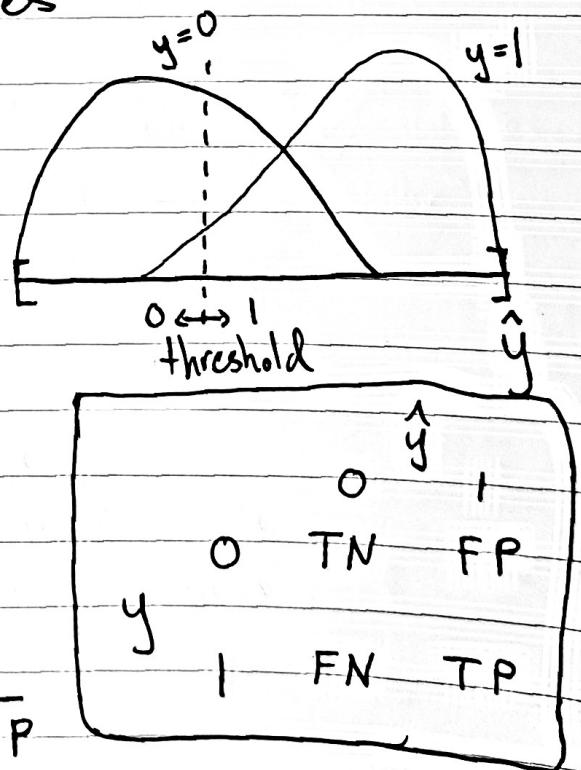
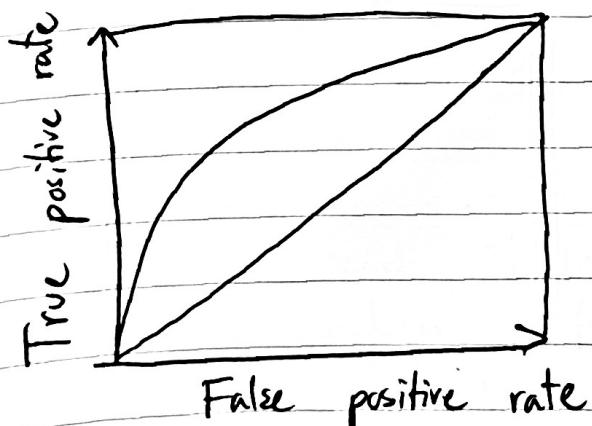
$$\text{min. } - \sum_n \ln \left[ \hat{y}_n^{y_n} \cdot (1 - \hat{y}_n)^{1-y_n} \right] + \underbrace{\|w\|_p^2}_{\text{regularization}}$$



$$\hat{y} = f(w_1 + w_2 x)$$

- logistic regression is classification,  
but provides probabilities

- Poisson ROC CURVE (metric)



$$TPR = \frac{TP}{TP+FN} \quad FPR = \frac{FP}{TN+FP}$$

Threshold @ 0  $\Rightarrow$  all 1  $\Rightarrow$  TPR=1, FPR=1  $\Rightarrow$  top right

## Sklearn pipelines

- Estimator (fit)
- Predictor (fit, predict)
- Transformer (fit, transform)

### Pipelines

$$T \rightarrow T \rightarrow T \rightarrow \dots \rightarrow T$$
$$T \rightarrow T \rightarrow T \rightarrow \dots \rightarrow T \rightarrow P$$

fit = fit, then transform, then next

transform = transform all

predict = transform all, predict

Feature Union can be used in parallel

### Example

SimpleImputer (strategy="mean")



OneHotEncoder()



PolynomialFeatures (degree=2)



Lasso (alpha=2)



Hyperparameters

- strategy
- degree
- alpha

Grid search

Random grid search

Bayesian Optimization

K-folds crossvalidation

## Tree models

- can be used for classification / regression
  - "20 questions" and splitting variables
  - very interpretable ("white box" model)
  - many hyperparameters
- 
- bagging : average over many trees
  - boosting : train tree n on the errors made by tree n-1.