

Primer on ML [1]

Overview

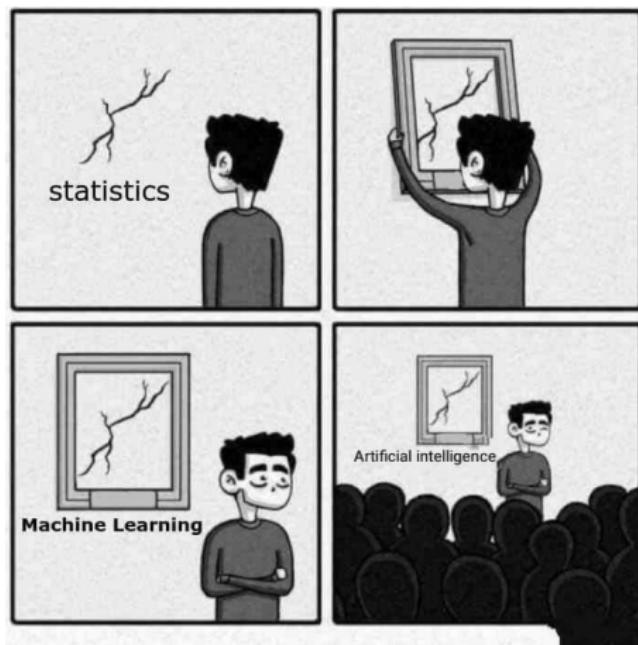
- What is ML
 - supervised [s]
 - unsupervised [u]
 - reinforcement learning [rl]
- Gently into [s]
 - linear models
 - trees
- ML specifics
 - metrics
 - overfitting
 - regularization
 - model validation

based on and inspired by

- MIPT x Yandex course [1]
- ML in HEP summer school [2]
- Data Mining in Action [3]

General concepts

What is ML

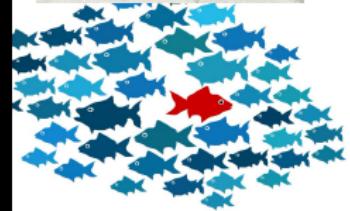
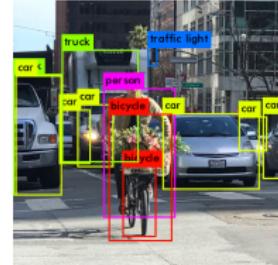


What is ML

Words

- Science of getting computers to act without being explicitly programmed [Coursera]
- Scientific study of algorithms and statistical models that computer systems use to effectively perform a specific task without using explicit instructions, relying on patterns and inference instead [Wikipedia]
- Science searching for hidden relations in data [A. Volokhova]

What is ML Pics



≡ Google Translate

<https://thiscatdoesnotexist.com/>

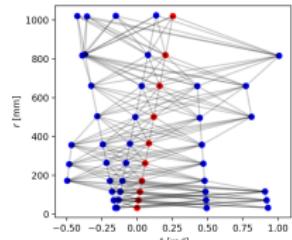
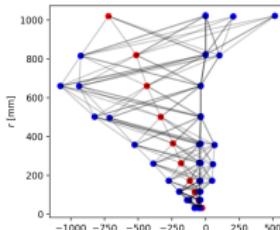
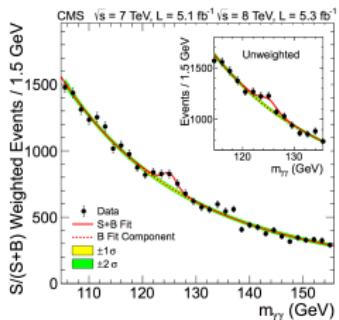
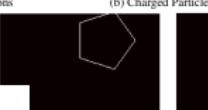
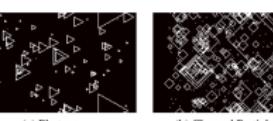
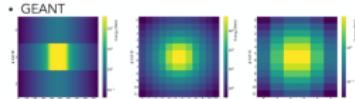
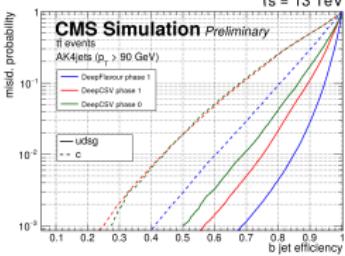
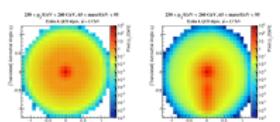
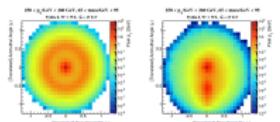
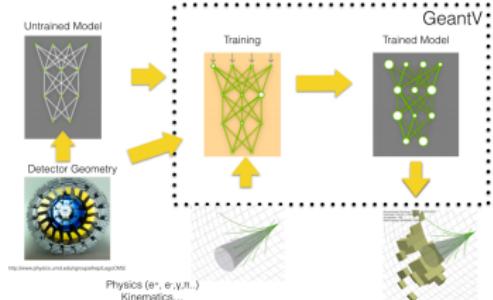
Text Documents

DETECT LANGUAGE RUSSIAN ENGLISH FRENCH ▾ + ENGLISH RUSSIAN SPANISH ▾

Translation

What is ML

HEP pics



- ▶ Supervised
 - Classification
 - Regression
- ▷ Unsupervised
 - Clustering
 - Anomaly detection
 - Generative models
 - Dimensionality reduction
 - ...
- ▷ Reinforcement Learning

Supervised Definitions

Feature space \mathbb{X}

- real
- categorical

Target space \mathbb{Y}

- $\mathbb{Y} = \overline{(0, C)}$ → classification
- $\mathbb{Y} = \mathbb{R}^d$ → regression

Sample $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell)$

Training

$$a(\mathbf{x}) = \operatorname{argmin}_{a \in \mathbb{A}} L(a, Y)$$

Training set $X, Y = (\mathbf{x}_i, \mathbf{y}_i)_{i=1}^\ell$

Model $a(\mathbf{x}) \leftarrow$ predicts target

Loss function $L(a, Y) \leftarrow$ evaluates model's quality

Supervised

C vs R

- ? cat, dog or muffin
 - rent price
 - temperature
 - annual profit
 - driving time
- ? relevant or spam
- ? disease or not
- ? good or bad

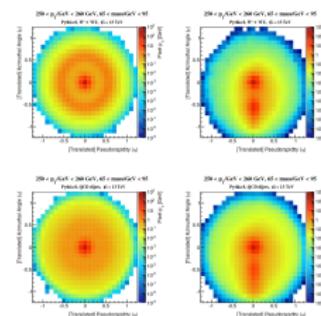
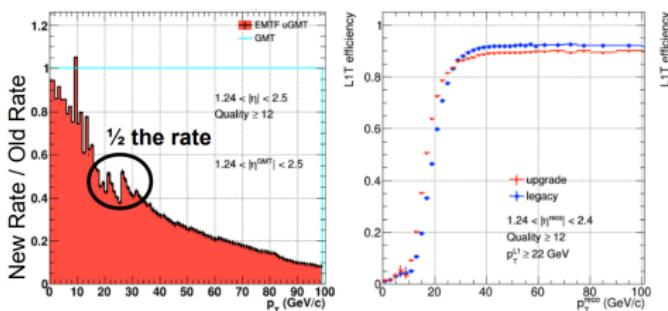
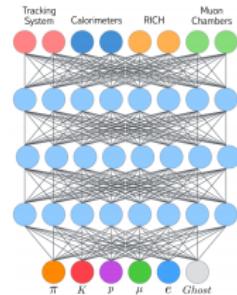
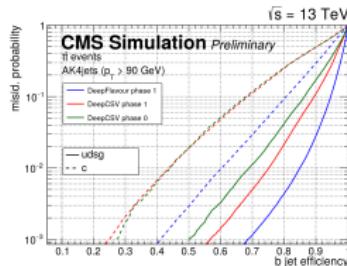
Supervised Classification in HEP

? b, c or ud jet

? π , K, or μ particle

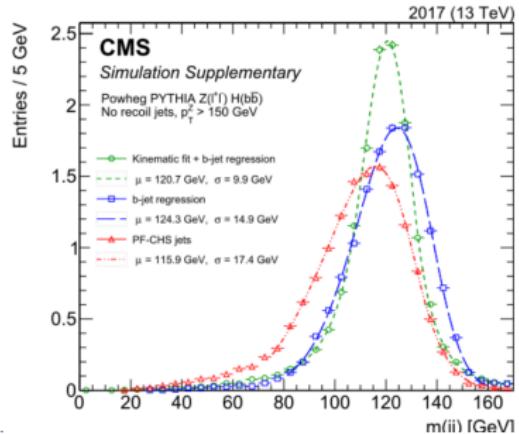
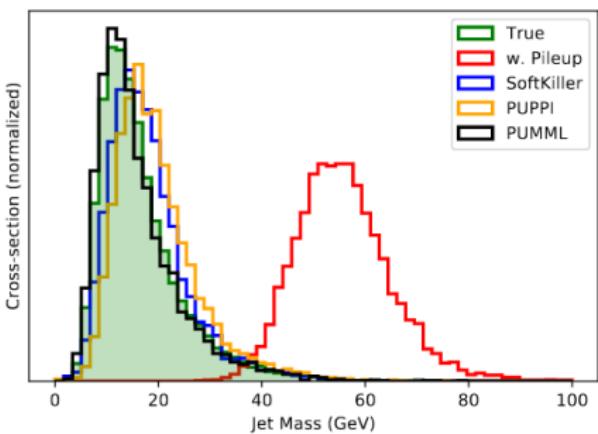
? $t\bar{t}$ or QCD event

? select or reject trigger candidate



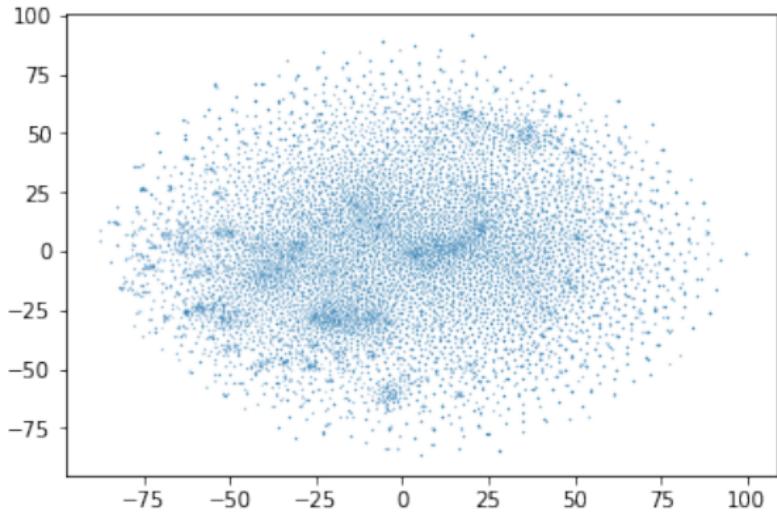
Supervised Regression in HEP

- energy resolution [1]
- PU mitigation [2]



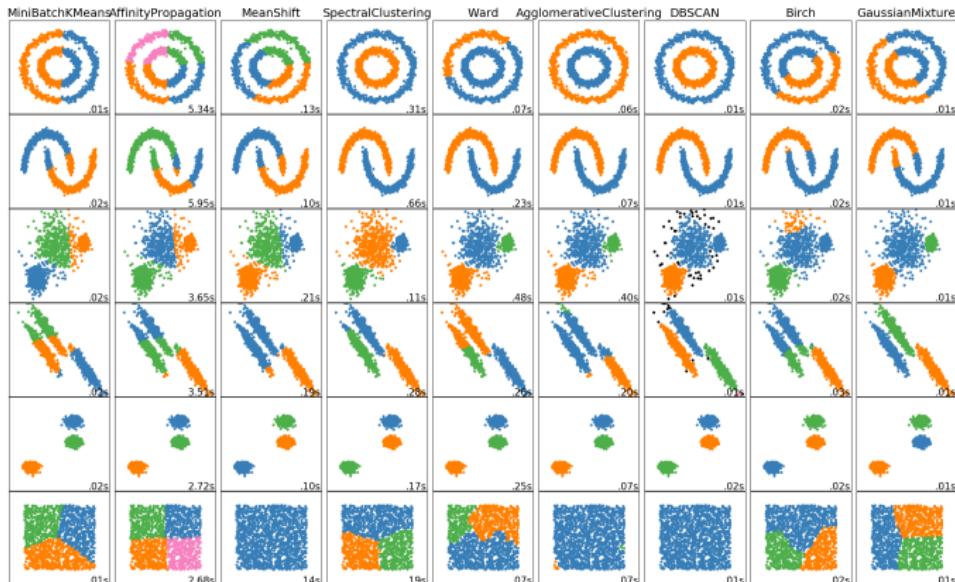
- ▷ **Supervised**
 - Classification
 - Regression
- ▶ **Unsupervised**
 - Clustering
 - Anomaly detection
 - Generative models
 - Dimensionality reduction
 - ...
- ▷ **Reinforcement Learning**

Unsupervised Clustering

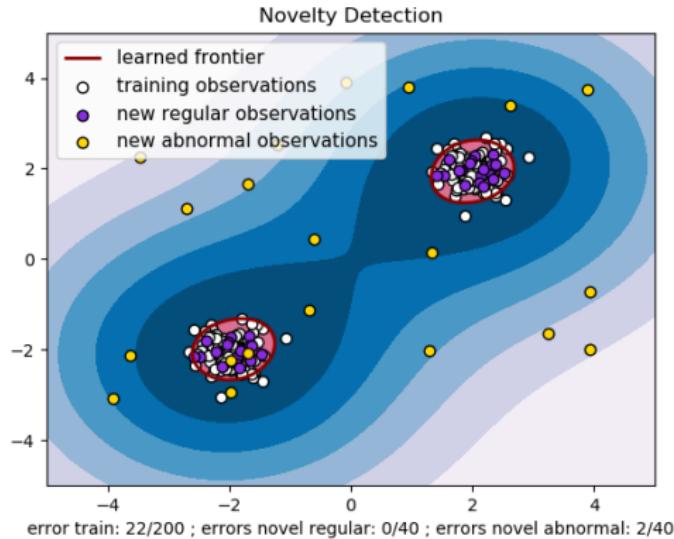


visuals! [link]

Unsupervised Clustering



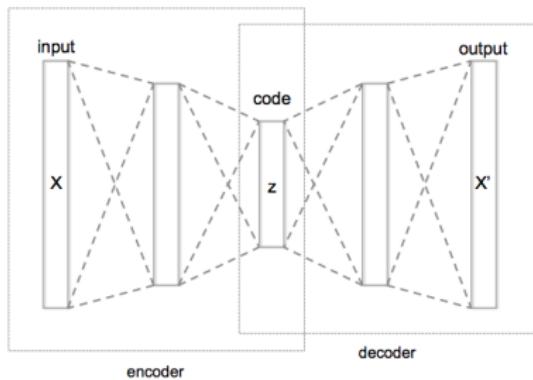
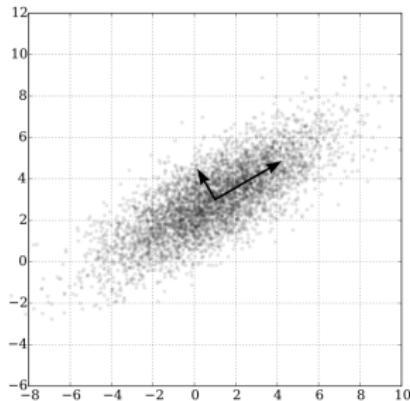
Unsupervised Novelty & Outliers



Unsupervised

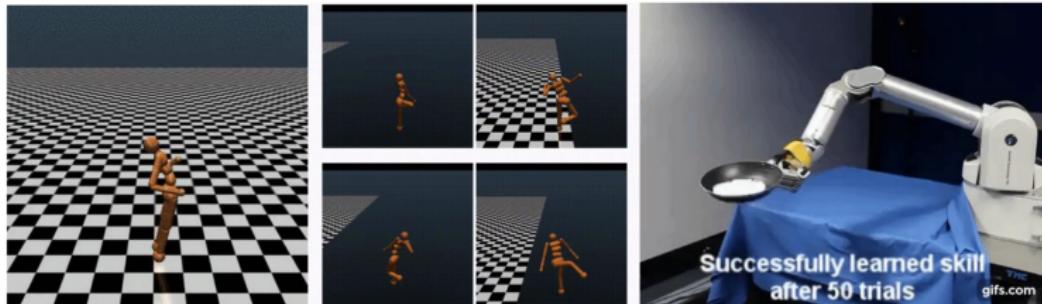
- Dimensionality reduction
 - Principal component analysis
 - Autoencoders
- * Generative models
 - GANs, VAE

more visuals! [\[link\]](#)

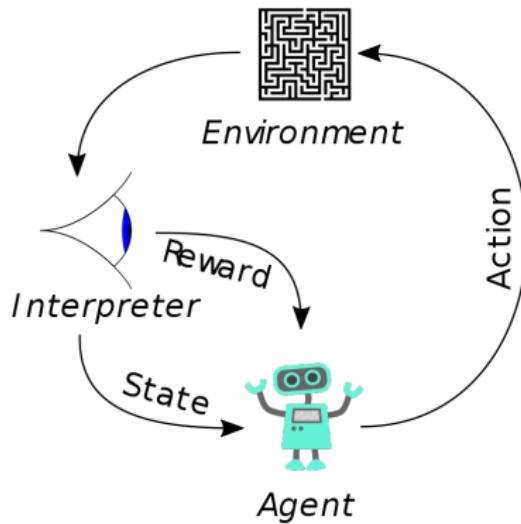


- ▷ **Supervised**
 - Classification
 - Regression
- ▷ **Unsupervised**
 - Clustering
 - Anomaly detection
 - Generative models
 - Dimensionality reduction
 - ...
- ▶ **Reinforcement Learning**

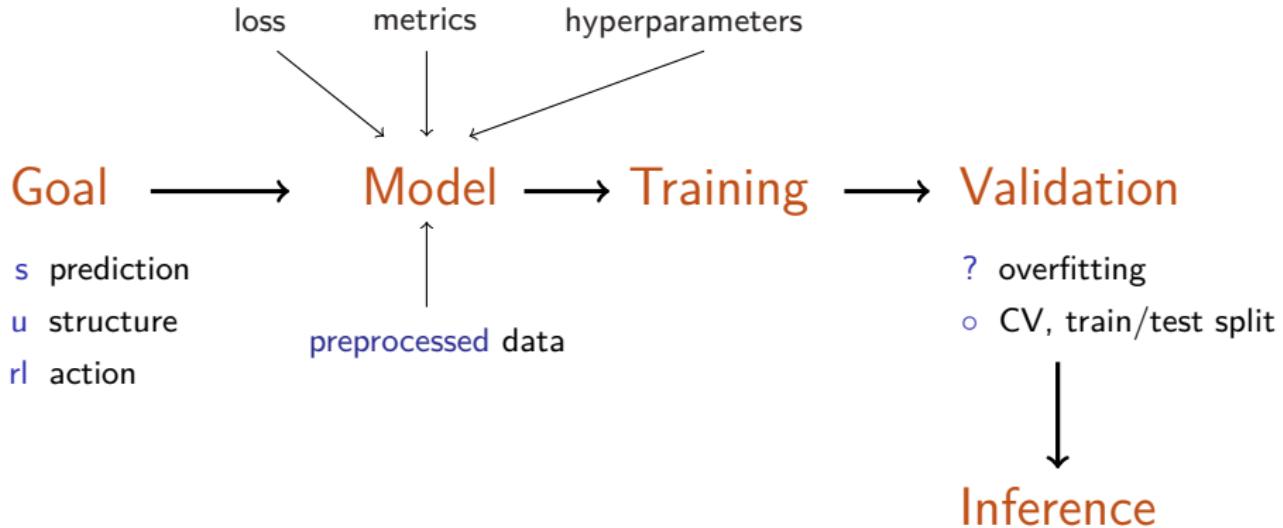
Reinforcement learning



Reinforcement learning



Pipeline

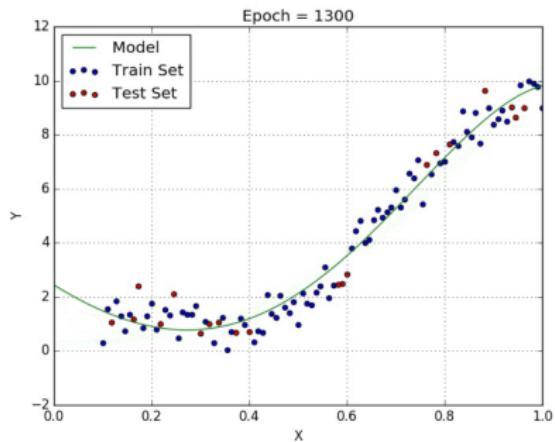


Linear models

Linear models

Regression

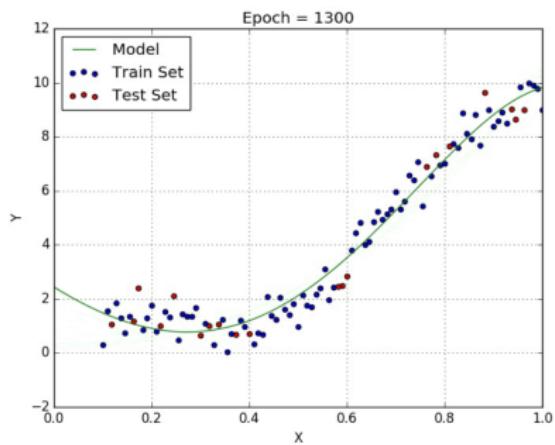
- $\mathbb{Y} = \mathbb{R}^d$
- $a(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$
- $L(a, Y)$:
 - MSE (L2) = $\frac{1}{\ell} \sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2$
 - RMSE = $\sqrt{\frac{1}{\ell} \sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2}$
 - MAE (L1) = $\frac{1}{\ell} \sum_{i=1}^{\ell} |y_i - a(\mathbf{x}_i)|$
 - $R^2 = 1 - \frac{\sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2}{\sum_{i=1}^{\ell} (y_i - \bar{y})^2}$
- $L(a, Y) \rightarrow \min_w$



Linear models

Regression

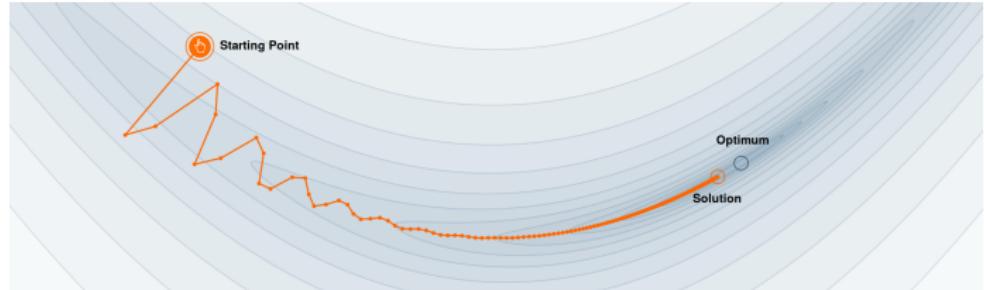
- $\mathbb{Y} = \mathbb{R}^d$
- $a(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$
- $L(a, Y)$:
 - MSE (L2) = $\frac{1}{\ell} \sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2$
 - RMSE = $\sqrt{\frac{1}{\ell} \sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2}$
 - MAE (L1) = $\frac{1}{\ell} \sum_{i=1}^{\ell} |y_i - a(\mathbf{x}_i)|$
 - $R^2 = 1 - \frac{\sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2}{\sum_{i=1}^{\ell} (y_i - \bar{y})^2}$



- $L(a, Y) \rightarrow \min_w \iff \text{how can we do that?}$

Linear models

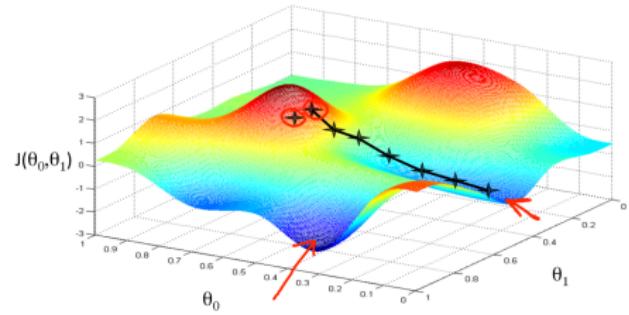
Gradient descent



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

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right)$$

! $\mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}^{(k-1)})$



more fancy [1], [2]

- $\mathbb{Y} = \{-1, 1\}$
- $a(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$
- $L(a, Y)$:
 - Accuracy = $\frac{1}{\ell} \sum_{i=1}^{\ell} [y_i = \text{sign}(\mathbf{w}^T \mathbf{x})]$

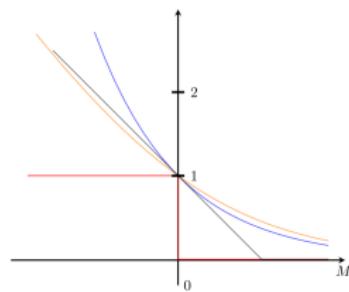
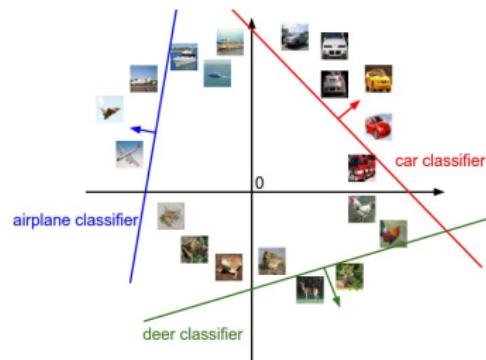
$$? \quad \frac{1}{\ell} \sum_{i=1}^{\ell} [y_i = a(\mathbf{x}_i)] \rightarrow \min_{\mathbf{w}}$$

! Optimize upper bound $\tilde{L}(a, Y)$

- Log loss
- Hinge loss
- ...

Linear models

Classification



and more fancy [1], [2]

Linear models

Logistic regression

$$\mathbb{Y} = \{-1, +1\}$$

- $a(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \equiv \sigma(\mathbf{w}^\top \mathbf{x})$

- $L(a, Y)$:

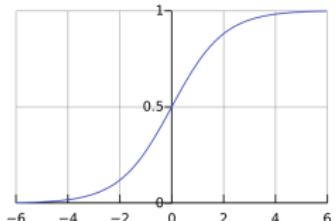
- Probability: $a(\mathbf{x}_i)^{[y_i=+1]} \cdot (1 - a(\mathbf{x}_i))^{[y_i=-1]}$

- Likelihood: $\mathcal{L} = \prod_{i=1}^{\ell} a(\mathbf{x}_i)^{[y_i=+1]} \cdot (1 - a(\mathbf{x}_i))^{[y_i=-1]}$

- Log likelihood: $\log \mathcal{L} = \sum_{i=1}^{\ell} [y_i = +1] \log a(\mathbf{x}_i) + [y_i = -1] \log(1 - a(\mathbf{x}_i))$

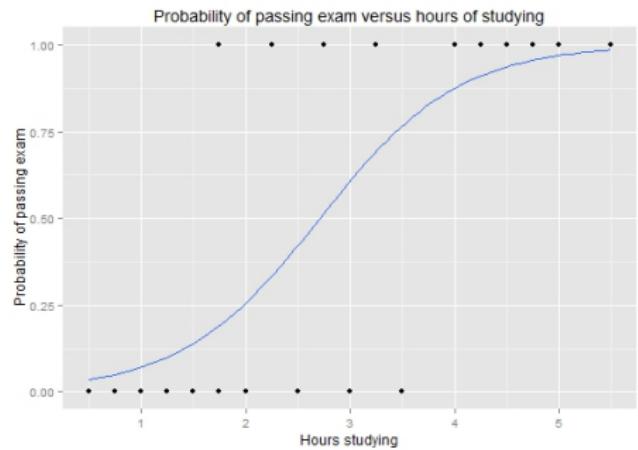
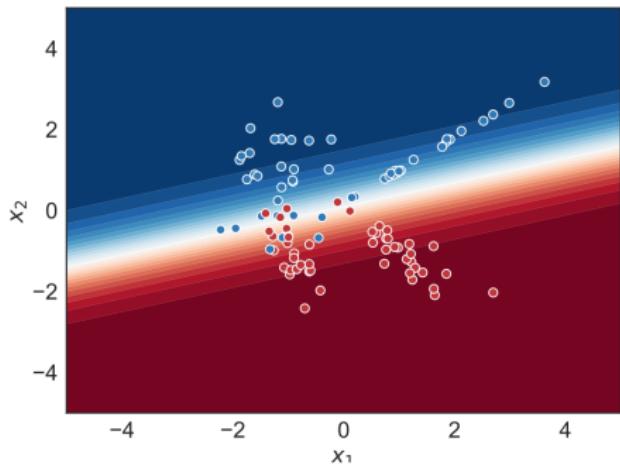
! predicts $P(y = +1 | \mathbf{x})$

! interpretable \mathbf{w} (in general for linear models)



Linear models

Logistic regression



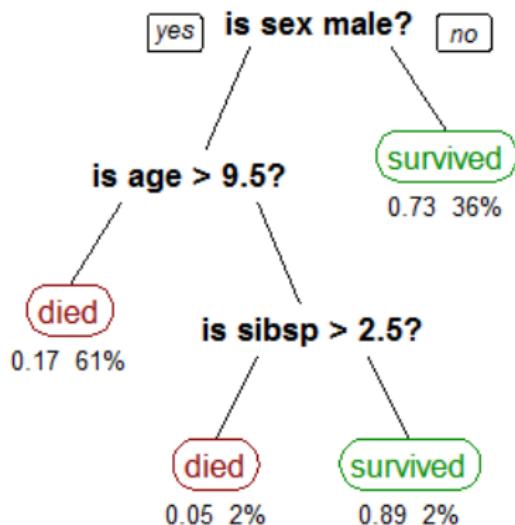
Linear models

pros & cons

- + Interpretable
- + Unique solution exists
- + Nice statistical properties
- + Fast
- + Perform well with lots of features
- + OK with little data
- Don't capture non-linear correlations
- Require feature linear independence

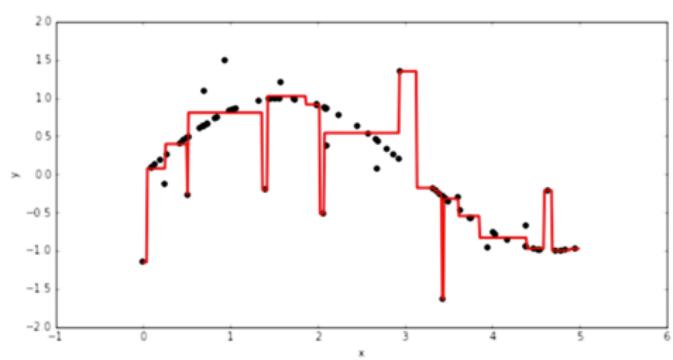
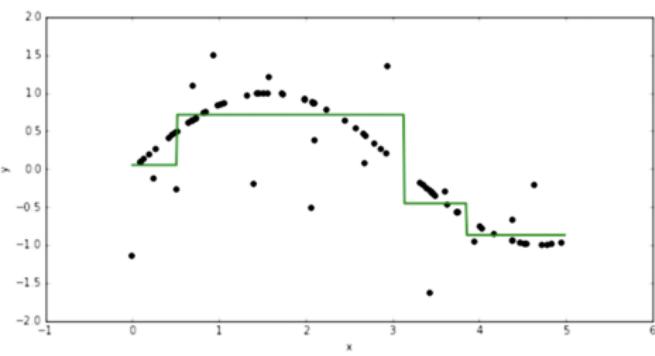
Trees

Trees



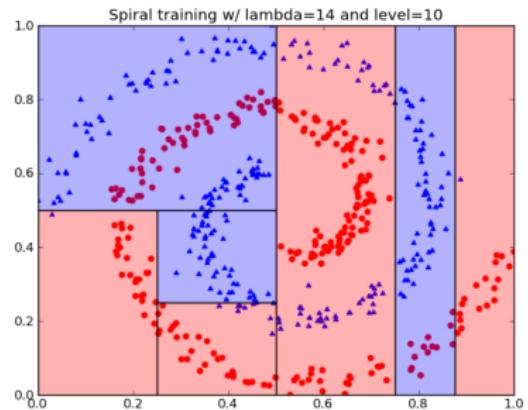
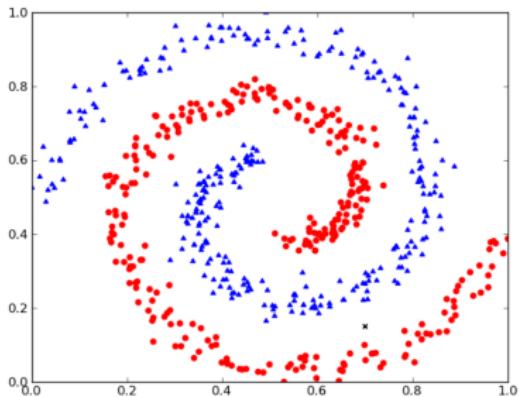
Trees

Regression



Trees

Classification

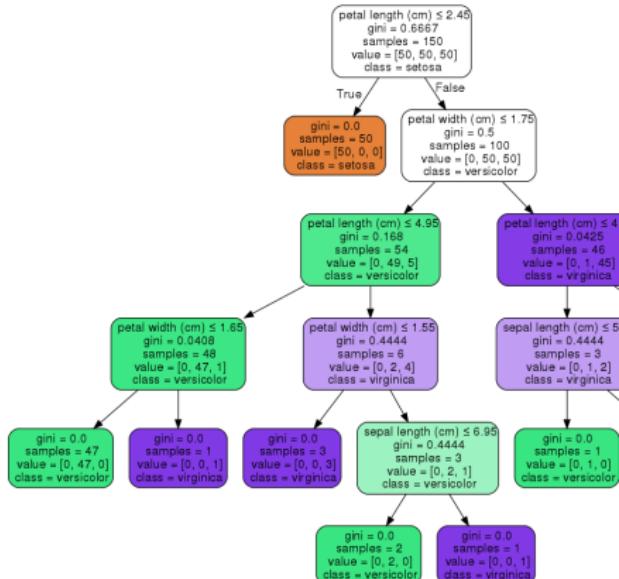


awesome visualisation [link]

Trees

Training

- Greedily split into Q_ℓ and Q_r :
 - $Q_\ell(\theta) = \{x, y \mid x_j \leq t_m; \theta = (j, t_m)\}$
 - $Q_r = Q \setminus Q_\ell(\theta)$
- $G(Q, \theta) = \frac{n_\ell}{N} H(Q_\ell(\theta)) + \frac{n_r}{N} H(Q_r(\theta))$
- Impurity $H(Q)$:
 - c Gini
 - c Entropy
 - r MSE (target variance), MAE
- minimize $\theta^* = \operatorname{argmin}_\theta G(Q, \theta)$
- stopping criteria
- leaf prediction - voting/mean value

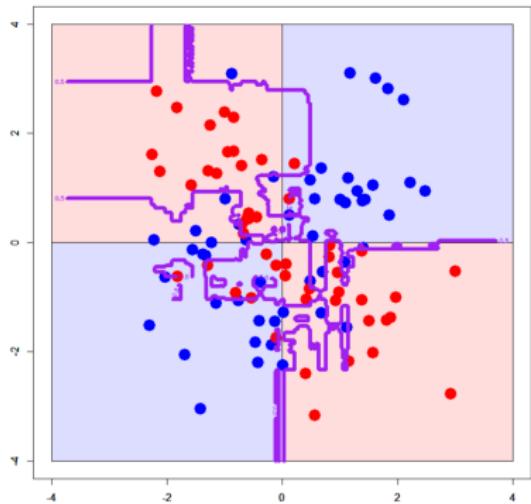


further reading [1], [2]

Trees

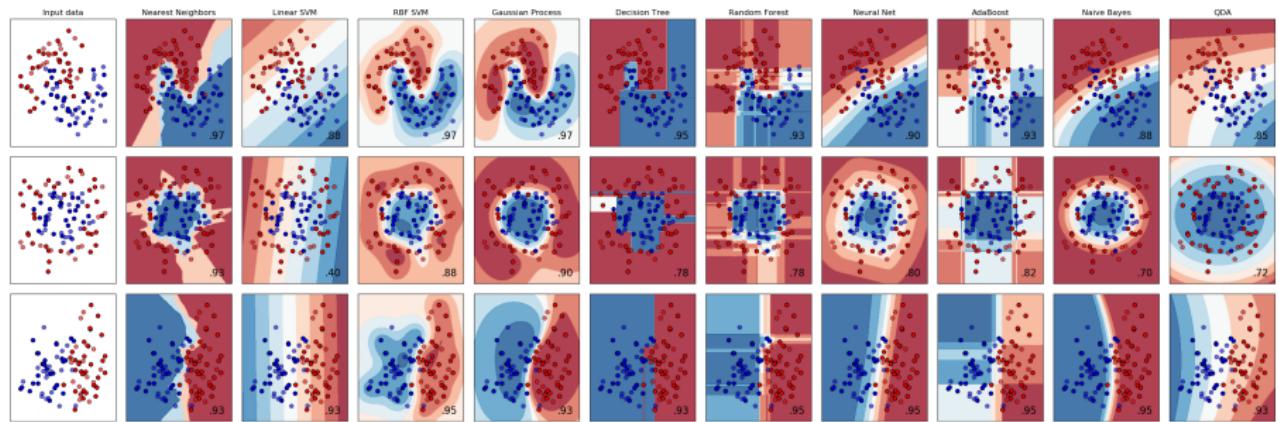
pros & cons

- + Interpretable
- + Categorical features
- + Feature linear dependence OK
- Overfits easily
- Not that powerful



comparison here [link]

there's more!



Out of the scope: Naive Bayes, kNN, SVM

ML specifics

They are **really** important

loss function

model validation

hyperparameter
optimization

nice one on metrics [1], [2]

- **MSE** = $\frac{1}{\ell} \sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2$
 - sensitive to outliers
 - easy to optimize
- **MAE** = $\frac{1}{\ell} \sum_{i=1}^{\ell} |y_i - a(\mathbf{x}_i)|$
 - robust to outliers
 - gradient problems (no)

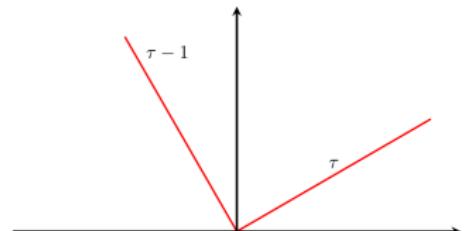
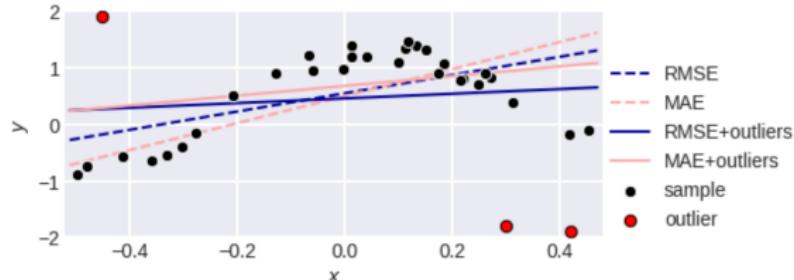
$$\bullet R^2 = 1 - \frac{\sum_{i=1}^{\ell} (y_i - a(\mathbf{x}_i))^2}{\sum_{i=1}^{\ell} (y_i - \bar{y})^2}$$

- explained variance
- comparison with naive mean model

- Logloss

Metrics

Regression



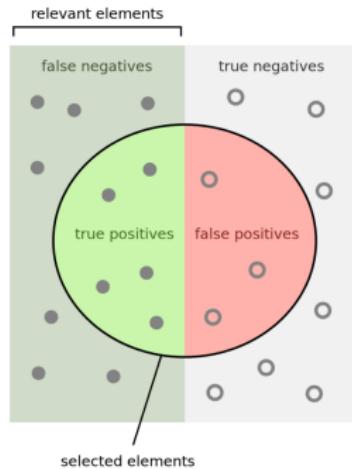
further reading [link]

- Accuracy = $\frac{1}{\ell} \sum_{i=1}^{\ell} [y_i = a(\mathbf{x}_i)]$
 - imbalanced data: 990b vs 10s
 - equal error weight
- Confusion matrix
 - + more info
- Precision, Recall
- F-score = $\frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$
- ! they are all threshold dependent
- Logloss

	$y = 1$	$y = -1$
$a(x) = 1$	True Positive (TP)	False Positive (FP)
$a(x) = -1$	False Negative (FN)	True Negative (TN)

Metrics

Classification



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

cool interactive demo [1]

more facts [2]

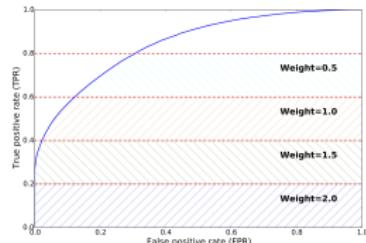
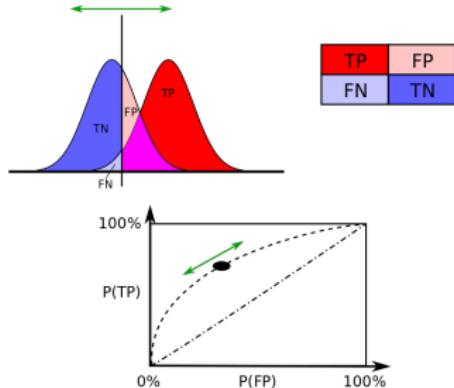
Metrics

ROC curve

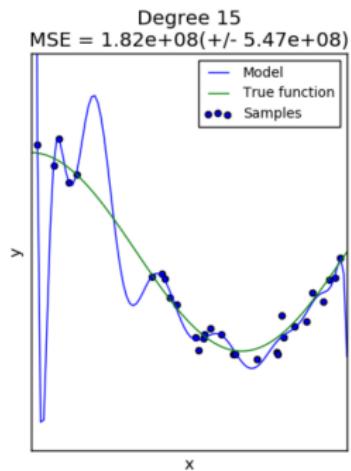
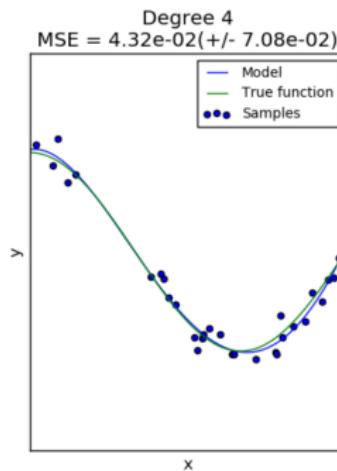
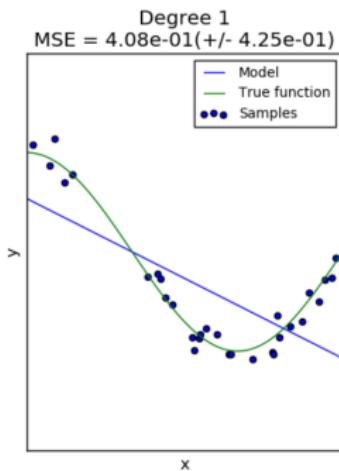
- curve in $(TPR(t), FPR(t))$ space
corr. to different thresholds $a(x) > t$

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

- ROC AUC = Area Under Curve
 - describe classifier as a whole
 - class separation (but not proba precision)
 - probability of correct pairs ordering
- weighted ROC AUC [3]

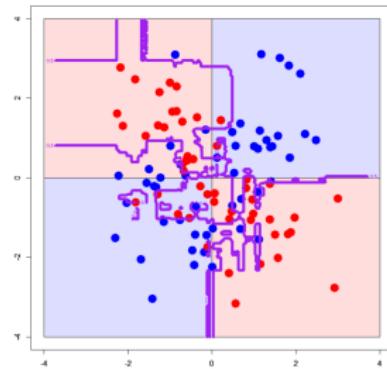


Overfitting



Overfitting

- Often looks like
 - $a(x) = 0.5 + 12458922x + 43983740x^2 + \dots + 2740x^9$
→ regularization
- Generalization:
 - underfitting - bad on train, bad on unseen
 - overfitting - excellent on train, bad on unseen
→ validate on hold-out data



Regularization

- Add penalty term to loss function:

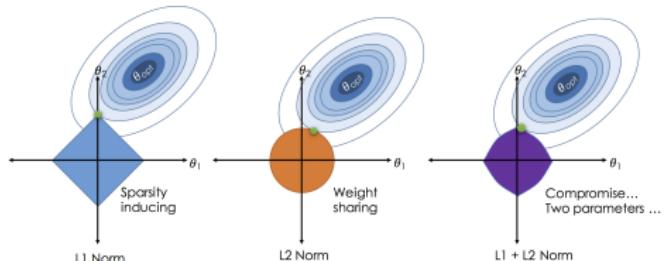
$$\text{L1: } \|w\|_1 = \sum_i |w_i|$$

$$\text{L2: } \|w\|_2 = \sum_i w_i^2$$

- $L(w, X) + \lambda \|w\|_{1,2} \rightarrow \min_w$

- λ is a **hyperparameter** of the model

! L1 zeros redundant w



Model validation

- ? check overfitting
- ? estimate model performance
- ? optimize hyperparameters

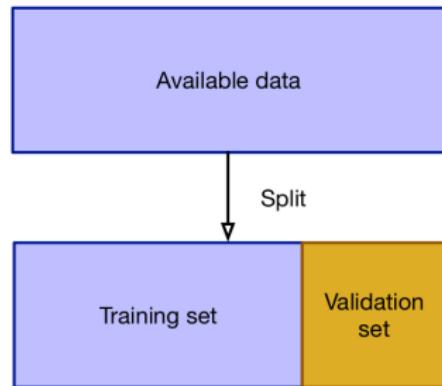
Solution

- o train/test split
- o cross validation (CV)

Model validation

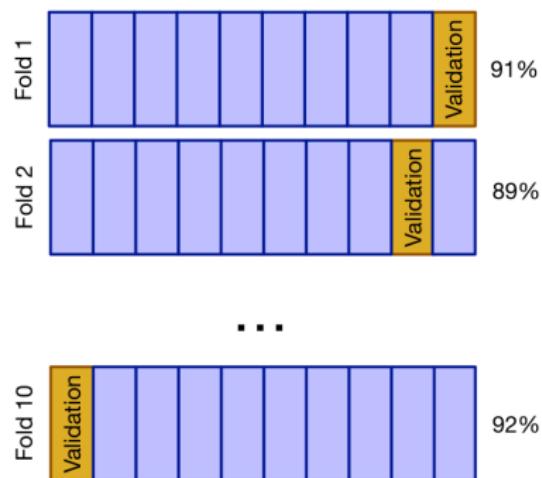
train/test

- might be luxury
- suffer from bad splitting
- **you** may overfit



- $CV = \frac{1}{K} \sum_{i=1}^K Q(a, X_k)$
 - + sensitive to overfitting
 - + can use all the data
 - still may suffer from bad splitting
 - many folds
 - + small bias
 - large variance
 - time-consuming
 - few folds
 - + small variance
 - + computationally effective
 - large bias
 - often used for hyperparameter optim.
 - e.g., polynomial degree d
 - train/test(CV)/val split

Model validation





Summary

- What is ML
 - supervised [s]
 - unsupervised [u]
 - reinforcement learning [rl]
- Gently into [s]
 - linear models
 - trees
- ML specifics
 - metrics
 - overfitting
 - regularization
 - model validation

next lecture

- Tree ensembles
- Feedforward NN
- * CNN & RNN