

### Machine Learning: The Probabilistic Perspective

#### Konstantin Tretyakov

http://kt.era.ee

IFI Summer School 2014 STACC 5

Software Technology and Applications Competence Center







### The Land of Machine Learning





### So far...

- Machine learning is important and interesting
- The general concept:

## Fitting models to data



### So far...

- Machine learning is important and interesting
- The general concept:





### So far...

- Instance-based methods
- Tree learning methods
- The "soul" of machine learning:

$$\operatorname{argmin}_{w} \operatorname{Error}(\operatorname{Data}, w) + \lambda \operatorname{Complexity}(w)$$

- Particular models:
  - OLS regression ( $\ell_2$ -loss, 0-penalty regression)
  - Ridge regression ( $\ell_2$ -loss,  $\ell_2$ -penalty regression)



#### Next

# Why should the model, tuned on the training set, generalize to the test set?



Learning purely from data is, in general, impossible

X	Υ	Output
0	0	False
0		True
I	0	True
1		? Free!
		IFI Summer School. June 2014



Learning purely from data is, in general, impossibleIs it good or bad?

What should we do to enable learning?



Learning purely from data is, in general, impossible

- Is it good or bad?
  - Good for cryptographers, bad for data miners
- What should we do to enable learning?
  - Introduce assumptions about data ("inductive bias"):
    - I. How does existing data relate to the future data?
    - 2. What is the system we are learning?



Learning purely from data is, in general, impossible

- Is it good or bad?
  - Good for cryptographers, bad for data miners
- What should we do to enable learning?
  - Introduce assumptions about data ("inductive bias").
    - How does existing data relate to the future data?
    - 2. What is the system we are learning?



#### How does existing data relate to future data?







heads, heads, tails, heads, tails,











### Probability theory



June 2014

-		
v	• T · E Probability distributions	[hide]
	Discrete univariate with finite support	[hide]
	Benford • Bernoulli • Beta-binomial • binomial • categorical • hypergeometric • Poisson binomial • Rademacher • discrete uniform • Zipf • Zipf-Mandelbro	ot
	Discrete univariate with infinite support	[hide]
	beta negative binomial • Boltzmann • Conway-Maxwell-Poisson • discrete phase-type • Delaporte • extended negative binomial • Gauss-Kuzmin • geome logarithmic • negative binomial • parabolic fractal • Poisson • Skellam • Yule-Simon • zeta	etric •
	Continuous univariate supported on a bounded interval, e.g. [0,1]	[hide]
	Arcsine • ARGUS • Balding-Nichols • Bates • Beta • Beta rectangular • Irwin-Hall • Kumaraswamy • logit-normal • Noncentral beta • raised cosine • triang U-quadratic • uniform • Wigner semicircle	ular •
	Continuous univariate supported on a semi-infinite interval, usually [0,∞)	[hide]
	Benini · Benktander 1st kind · Benktander 2nd kind · Beta prime · Bose-Einstein · Burr · chi-squared · chi · Coxian · Dagum · Davis · Erlang · exponential Fermi-Dirac · folded normal · Fréchet · Gamma · generalized inverse Gaussian · half-logistic · half-normal · Hotelling's T-squared · hyper-exponential hypoexponential · inverse chi-squared (scaled-inverse-chi-squared) · inverse Gaussian · inverse gamma · Kolmogorov · Lévy · log-Cauchy · log-Lapla log-logistic · log-normal · Maxwell-Boltzmann · Maxwell speed · Mittag-Leffler · Nakagami · noncentral chi-squared · Pareto · phase-type · Rayleigh relativistic Breit-Wigner · Rice · Rosin-Rammler · shifted Gompertz · truncated normal · type-2 Gumbel · Weibull · Wilks' lambda	· F ·
	Continuous univariate supported on the whole real line (-∞, ∞)	[hide]
Cau L	uchy • exponential power • Fisher's z • generalized normal • generalized hyperbolic • geometric stable • Gumbel • Holtsmark • hyperbolic secant • Landau • .innik • logistic • noncentral t • normal (Gaussian) • normal-inverse Gaussian • skew normal • slash • stable • Student's t • type-1 Gumbel • variance-gamma	Laplace • Voigt
	Continuous univariate with support whose type varies	[hide]
	generalized extreme value · generalized Pareto · Tukey lambda · q-Gaussian · q-exponential · shifted log-logistic	
	Mixed continuous-discrete univariate distributions	[hide]
	rectified Gaussian	
	Multivariate (joint)	[hide]
	Discrete: Ewens • multinomial • Dirichlet-multinomial • negative multinomial Continuous: Dirichlet • Generalized Dirichlet • multivariate normal • Multivariate stable • multivariate Student • normal-scaled inverse gamma • normal-gam Matrix-valued: inverse matrix gamma • inverse-Wishart • matrix normal • matrix t • matrix gamma • normal-inverse-Wishart • matrix • Wishar	nma t
	Directional	[hide]
	Univariate (circular) directional: Circular uniform • univariate von Mises • wrapped normal • wrapped Cauchy • wrapped exponential • wrapped Lév Bivariate (spherical): Kent • Bivariate (toroidal): bivariate von Mises Multivariate: von Mises–Fisher • Bingham	у
	Degenerate and singular	[hide]
	Degenerate: discrete degenerate · Dirac delta function Singular: Cantor	
	Families	[hide]
	June 2014	

F



### Probability theory

D





from numpy.random import beta, binomial, chisquare, dirichlet, exponential, f, gamma, geometric, gumbel, hypergeometric, ...

>>> numpy.random.seed(1)
>>> binomial(10, 0.2)
:::: 2



```
>>> numpy.random.seed(1)
>>> X = binom(10, 0.2)
>>> X.rvs()
::: 2
```

>>> X.pmf(2), X.cdf(2), X.mean(), X.std(), ...



#### What is your height?





What is your height?

#### Is it a fixed number?





### What is your height?

### Is it a fixed number?

- Frequentist: Yes, it is, we just don't know it precisely.
- Bayesian: No, it is not. It is a distribution.





### What is your height?

### Is it a fixed number?

- Frequentist: Yes, it is, we just don't know it precisely.
- Bayesian: No, it is not. It is a distribution.

In any case, we need probabilistic reasoning.





#### Statistics

How do we **infer** a probabilistic **model** based on data?







#### Statistics

How do we **infer** a probabilistic **model** based on data?





#### Statistics

How do we **infer** a probabilistic **model** based on data?







#### Decision theory

How do we **use** a probabilistic model **to predict**?





#### Decision theory

How do we **use** a probabilistic model **to predict**?







Model, trained on the training set might work well on the test set because:

Because we **assume** a single underlying mechanism.

Because we use statistical inference to infer the mechanism.

Because we use decision theory to produce optimal decisions.

### Quiz





#### **Statistics**









IFI Summer School. June 2014





#### Space of candidate models







#### **Statistics**













#### Hypothesis testing







or not?





#### **Model selection**











#### **Parameter inference**







### Parameter inference



#### **Biased coin**




## Data Likelihood: Pr[Data | Model]

#### Example:

- Model: Be(0.5)
- Data: 1,1,0,1,1

Likelihood: ?



## Data Likelihood: Pr[Data | Model]

#### Example:

- Model: Be(0.5)
- Data: 1,1,0,1,1
- Likelihood:  $0.5 \cdot 0.5 \cdot 0.5 \cdot 0.5 \cdot 0.5 = 2^{-5}$

0.03125



## Data Likelihood: Pr[Data | Model]

#### • Example:

- Model: Be(0.2)
- Data: 1,1,0,1,1

Likelihood: ?



## Data Likelihood: Pr[Data | Model]

#### • Example:

- Model: Be(0.2)
- Data: 1,1,0,1,1
- Likelihood:  $0.2 \cdot 0.2 \cdot 0.8 \cdot 0.2 \cdot 0.2 = 0.2^4 \cdot 0.8$

0.00128



#### • Example:

- Model: Be(p)
- Data: 1,1,0,1,1
- Likelihood:  $p \cdot p \cdot (1-p) \cdot p \cdot p = p^{n_1}(1-p)^{n_0}$





#### • Example:

- Model: Be(p)
- Data: 1,1,0,1,1
- Likelihood:  $p \cdot p \cdot (1-p) \cdot p \cdot p = p^{n_1}(1-p)^{n_0}$





## argmax<sub>Model</sub> Pr(Data |Model)



- You are on a trip in an exotic country and you meet a person who happens to be from Switzerland.
- Is he a member of the Swiss Parliament?



- Data: "X is from Switzerland"
- Models:
  - "X is a member of Swiss Parliament",
  - "X is not a member of Swiss Parliament"



- Data: "X is from Switzerland"
- Models:
  - "X is a member of Swiss Parliament",
  - "X is not a member of Swiss Parliament"

#### Likelihoods:

- P(X is from Switzerland | X is a member of SP) =
- P(X is from Switzerland | X is not a member of SP) =



- Data: "X is from Switzerland"
- Models:
  - "X is a member of Swiss Parliament",
  - "X is not a member of Swiss Parliament"

#### Likelihoods:

- P(X is from Switzerland | X is a member of SP) = 1
- P(X is from Switzerland | X is **not** a member of SP) =  $\frac{8}{7000}$



- Data: "X is from Switzerland"
- Models:
  - "X is a member of Swiss Parliament",
- MLE treats all candidate models
  L as equal and can thus overfit
  P(X is from Switzerland | X is a member of SP) = 1
  - P(X is from Switzerland | X is **not** a member of SP) =  $\frac{8}{7000}$



Maximum Likelihood Estimate (MLE):

argmax<sub>Model</sub> Pr(Data | Model)

Maximum A-posteriori Estimate (MAP):

argmax<sub>Model</sub> Pr(||Data)



#### argmax<sub>Model</sub> Pr(Model|Data)



#### argmax<sub>Model</sub> Pr(Model|Data)

## $argmax_{Model} \frac{Pr(Model, Data)}{Pr(Data)}$

#### argmax<sub>Model</sub> Pr(Model, Data)



D

#### argmax<sub>Model</sub> Pr(Model|Data)

## $argmax_{Model} \frac{Pr(Model, Data)}{Pr(Data)}$

#### argmax<sub>Model</sub> Pr(Model, Data)

#### argmax<sub>Model</sub> Pr(Data | Model) · Pr(Model)



## MAP Estimation argmax<sub>Model</sub> Pr(Model|Data) argmax<sub>Model</sub> <u>Pr(Model</u>, <u>Model posterior</u> Pr(Data)

#### argmax<sub>Model</sub> Pr(Model, Data)





#### Summary

Maximum Likelihood Estimate (MLE):

argmax<sub>Model</sub> Pr(Data | Model)

Maximum A-posteriori Estimate (MAP):

argmax<sub>Model</sub> Pr(Data | Odel) Pr(Model)



## Model: Be(p)

D

#### Data: 1,1,0,1,1





June 2014

RJU ÜLIK,



D







#### argmax<sub>Model</sub> Pr(Data | Model) · Pr(Model)



D

#### argmax<sub>Model</sub> Pr(Data | Model) · Pr(Model)

#### argmax<sub>Model</sub> log (Pr(Data | Model) · Pr(Model))



D

#### argmax<sub>Model</sub> Pr(Data | Model) · Pr(Model)

#### argmax<sub>Model</sub> log (Pr(Data | Model) · Pr(Model))

#### $\operatorname{argmax}_{\operatorname{Model}} \log \Pr(\operatorname{Data}|\operatorname{Model}) + \log \Pr(\operatorname{Model})$



argmax<sub>Model</sub> Pr(Data | Model) · Pr(Model)

argmax<sub>Model</sub> log (Pr(Data | Model) · Pr(Model))

argmax<sub>Mode</sub> log Pr(Data|Model) + log Pr(Model)



argmax<sub>Model</sub> Pr(Data | Model) · Pr(Model)

argmax<sub>Model</sub> log (Pr(Data | Model) · Pr(Model))

argmax<sub>Mode</sub> log Pr(Data|Model) + log Pr(Model)

argmin<sub>w</sub> Error(Data, w) + Complexity(w)



#### Problems of MAP estimation

Þ







#### Problems of MAP estimation



D



## Pick the model with minimal expected risk E(Model |Data)



D



## Pick the model with minimal expected risk E(Model | Data)





#### Bayesian estimation +

D

# Use the full posterior distribution Pr(Model |Data)





#### Quiz

D

#### Three major model inference methods are:



#### Linear Regression (again)

D

#### **Normal distribution**





#### Linear Regression (again)






#### $Y = \boldsymbol{w}^T \boldsymbol{x} + N(0, \sigma^2)$



 $Y = \boldsymbol{w}^T \boldsymbol{x} + N(0, \sigma^2)$ 

 $e = Y - \boldsymbol{w}^T \boldsymbol{x} \sim N(0, \sigma^2)$ 



 $Y = \boldsymbol{w}^T \boldsymbol{x} + N(0, \sigma^2)$  $e = Y - \boldsymbol{w}^T \boldsymbol{x} \sim N(0, \sigma^2)$  $\Pr((x, y) | \boldsymbol{w}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\frac{e^2}{\sigma^2}\right)$ 



 $Pr((x_1, y_1), (x_2, y_2), ..., (x_n, y_n) | \mathbf{w}, \sigma^2)$  $= \prod_{i} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$ 



 $Pr((x_1, y_1), (x_2, y_2), ..., (x_n, y_n) | \mathbf{w}, \sigma^2)$  $\propto \prod \exp\left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$ 

D



log Pr(( $x_1, y_1$ ), ( $x_2, y_2$ ), ..., ( $x_n, y_n$ )| $w, \sigma^2$ )  $\propto \log \prod \exp \left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$ 



log Pr( $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) | \mathbf{w}, \sigma^2$ )  $\propto \log \prod \exp \left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$  $=\sum \left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$ 



log Pr( $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) | \mathbf{w}, \sigma^2$ )  $\propto \log \prod \exp \left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$  $\left(-\frac{1}{2}\frac{e_i^2}{\sigma^2}\right)$  $=-\frac{1}{2\sigma^2}$ 



#### MLE and OLS

D

## $\operatorname{argmax}_{\mathbf{w}} \operatorname{LogLikelihood}(\operatorname{Data}, \mathbf{w}) = \operatorname{argmin}_{\mathbf{w}} \sum_{i} e_{i}^{2}$



#### $\Pr(w|\text{Data}) \propto \Pr(\text{Data}|w) \cdot \Pr(w)$



D

#### $\log \Pr(w | \text{Data}) \propto \log \Pr(\text{Data} | w) + \log \Pr(w)$



## $\log \Pr(w|\text{Data}) \propto \log \Pr(\text{Data}|w) + \log \Pr(w)$





## $\log \Pr(w|\text{Data}) \propto \log \Pr(\text{Data}|w) + \log \Pr(w)$



### Let the prior on $w_j$ be Gaussian: $Pr(w_j) \propto exp\left(-\frac{w_j^2}{2\alpha^2}\right)$



## $\log \Pr(w|\text{Data}) \propto \log \Pr(\text{Data}|w) + \log \Pr(w)$



### Let the prior on $w_j$ be Gaussian: $Pr(w) \propto \prod_{i} \exp\left(-\frac{w_j^2}{2\alpha^2}\right)$



## $\log \Pr(w|\text{Data}) \propto \log \Pr(\text{Data}|w) + \log \Pr(w)$



Let the prior on  $w_j$  be Gaussian:  $\log \Pr(w) \propto \sum_j -\frac{w_j^2}{2\alpha^2}$ 



## $\log \Pr(w|\text{Data}) \propto \log \Pr(\text{Data}|w) + \log \Pr(w)$



### Let the prior on $w_j$ be Gaussian: $\log \Pr(w) \propto -\sum_i w_j^2$





Let the prior on  $w_j$  be Gaussian:  $\log \Pr(w) \propto -\sum_i w_j^2$ 

























### Probability for modeling

#### Statistics for estimation



#### Decision theory for prediction





Þ



#### Decision





#### Decision





VALID

0

#### Decision

SPAM

0

5





VALID

0

**0.8** 

#### Decision

SPAM

0

5





VALID

0

**0.8** 

#### Decision





# $R(\hat{y}|x) = \sum_{y} \Pr(y|x)\ell(\hat{y},y)$

D









$$R(\hat{y}|x) = \int_{y} \ell(\hat{y}, y) dF(y|x)$$



#### **Bayesian Classifier**

• Optimal classifier:

For a given x and a conditional probabilistic model Pr(y|x)predict  $\hat{y}$ , that has the **smallest expected risk**.



### **Bayesian** Classifier

• Optimal classifier:

For a given x and a conditional probabilistic model Pr(y|x)predict  $\hat{y}$ , that has the **smallest expected risk**.

For symmetric risk  $\ell$ , this corresponds to picking the option with the highest probability.



### Probability for modeling

#### Statistics for estimation

### MLE MAP

#### Decision theory for prediction

#### **Bayesian Classifier**




#### The Tennis Dataset

Þ

Day	Outlook	Temp	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



#### Shall we play tennis today?

D

PlayTennis
No
No
Yes
Yes
Yes
No
Yes
No
Yes
No



#### Shall we play tennis today?

	Estimate a
PlayTennis	
No	probabilistic model
No	and predict:
Yes	
Yes	
Yes	$Pr(Y_{es}) = 9/14 = 0.64$
No	11(103) = 7711 = 0.01
Yes	Pr(No) = 5/14 = 0.36
No	
Yes	
Yes	Tes
Yes	
Yes	
Yes	
No	



Wind	PlayTennis
Weak	No
Strong	No
Weak	Yes
Weak	Yes
Weak	Yes
Strong	No
Strong	Yes
Weak	No
Weak	Yes
Weak	Yes
Strong	Yes
Strong	Yes
Weak	Yes
Strong	No



Wind	PlayTennis
Weak	No
<i>w</i> еак	NO
Strong	No
Weak	Yes
Weak	Yes
Weak	Yes
Strong	No
Strong	Yes
Weak	No
Weak	Yes
Weak	Yes
Strong	Yes
Strong	Yes
Weak	Yes
Strong	No

Pr(Yes | Weak) = 6/8 Pr(No | Weak) = 2/8

Pr(Yes | Strong) = 3/6Pr(No | Strong) = 3/6



Humidity	Wind	PlayTennis
High	Weak	No
High	Strong	No
High	Weak	Yes
High	Weak	Yes
Normal	Weak	Yes
Normal	Strong	No
Normal	Strong	Yes
High	Weak	No
Normal	Weak	Yes
Normal	Weak	Yes
Normal	Strong	Yes
High	Strong	Yes
Normal	Weak	Yes
High	Strong	No

Pr(Yes | High,Weak) = 2/4 Pr(No | High,Weak) = 2/4

Pr(Yes | High,Strong) = 1/3 Pr(No | High,Strong) = 2/3



#### The Bayesian Classifier

In general:

- **I. Estimate from data:**  $Pr(Class | x_1, x_2, x_3, ...)$
- 2. For a given instance  $(x_1, x_2, x_3, ...)$ predict class whose conditional probability is greater:

 $Pr(C_1 | x_1, x_2, x_3, ...) > Pr(C_2 | x_1, x_2, x_3, ...)$ → predict C<sub>1</sub>

#### Problem



#### • We need exponential amount of data

Humidity	Wind	PlayTennis
High	Weak	No
High	Strong	No
High	Weak	Yes
High	Weak	Yes
Normal	Weak	Yes
Normal	Strong	No
Normal	Strong	Yes
High	Weak	No
Normal	Weak	Yes
Normal	Weak	Yes
Normal	Strong	Yes
High	Strong	Yes
Normal	Weak	Yes
High	Strong	No

Pr(Yes | High,Weak) = 2/4 Pr(No | High,Weak) = 2/4

Pr(Yes | High,Strong) = 1/3 Pr(No | High,Strong) = 2/3



To scale beyond 2-3 attributes, use a hack:

## Assume that attributes are independent within each class:

#### $Pr(x_1, x_2, x_3 | Class)$ = $Pr(x_1 | Class) Pr(x_2 | Class) Pr(x_3 | Class) ...$



1.  $Pr(C_1|x) > Pr(C_2|x)$  $\rightarrow$  predict C<sub>1</sub>  $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$ 







1.  $Pr(C_1 | x) > Pr(C_2 | x)$  $\rightarrow$  predict C<sub>1</sub> 2.  $\frac{\Pr(C_1)\Pr(\boldsymbol{x}|C_1)}{\Pr(\boldsymbol{x})} > \frac{\Pr(C_2)\Pr(\boldsymbol{x}|C_2)}{\Pr(\boldsymbol{x})}$  $\rightarrow$  predict C<sub>1</sub> 3.  $Pr(C_1) Pr(x|C_1) > Pr(C_2) Pr(x|C_2)$  $\rightarrow$  predict C<sub>1</sub>



1.  $Pr(C_1 | x)$  $> \Pr(C_2|x)$  $\rightarrow$  predict C<sub>1</sub> 2.  $\frac{\Pr(C_1)\Pr(\boldsymbol{x}|C_1)}{\Pr(\boldsymbol{x})} > \frac{\Pr(C_2)\Pr(\boldsymbol{x}|C_2)}{\Pr(\boldsymbol{x})}$  $\rightarrow$  predict  $C_1$ 3.  $Pr(C_1) Pr(x|C_1) > Pr(C_2) Pr(x|C_2)$  $\rightarrow$  predict C<sub>1</sub> 4.  $Pr(C_1) \cdot Pr(x_1|C_1) Pr(x_2|C_1) \dots Pr(x_m|C_1) >$  $Pr(C_2) \cdot Pr(x_1|C_2) Pr(x_2|C_2) \dots Pr(x_m|C_2)$  $\rightarrow$  predict C<sub>1</sub>







#### Works for both discrete and continuous attributes.

#### The goods:

- Easy to implement, efficient
- Won't overfit, intepretable
- Works better than you would expect (e.g. spam filtering)

#### The bads

- "Naïve", linear
- Usually won't work well for too many classes
- Not a good probability estimator



# from sklearn.naive\_bayes import BernoulliNB, MultinomialNB, GaussianNB



MLE:

Quiz



MAP:

argmax<sub>Model</sub>\_\_\_\_\_

Gaussian distribution:  $f(x) = \text{const} \times \exp(\underline{\qquad})$ 



#### Bayesian classifier has optimal

111Z

#### Naïve Bayesian classifier assumption:

- $Pr(C|x_1, x_2) = Pr(C|x_1) Pr(C|x_2)$
- $\Pr(x_1, x_2 | C) = \Pr(x_1 | C) \Pr(x_2 | C)$
- $Pr(C_1, C_2|x) = Pr(C_1|x) Pr(C_2|x)$
- $\Pr(x|C_1, C_2) = \Pr(x|C_1) \Pr(x|C_2)$



#### All machine learning methods we have considered so far rely on MLE or MAP

• Yes

No



#### The Land of Machine Learning





### **Questions**?



THEN I TOOK A STATISTICS CLASS. NOW I DON'T.



http://xkcd.com/552/