

Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers

Meelis Kull, Telmo Silva Filho, Peter Flach

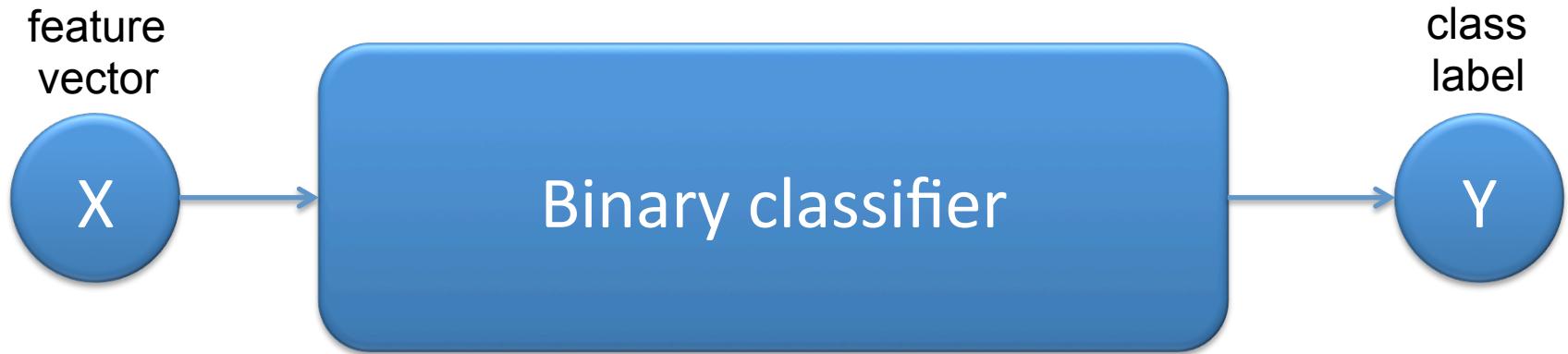
AISTATS 2017



www.irc-sphere.ac.uk

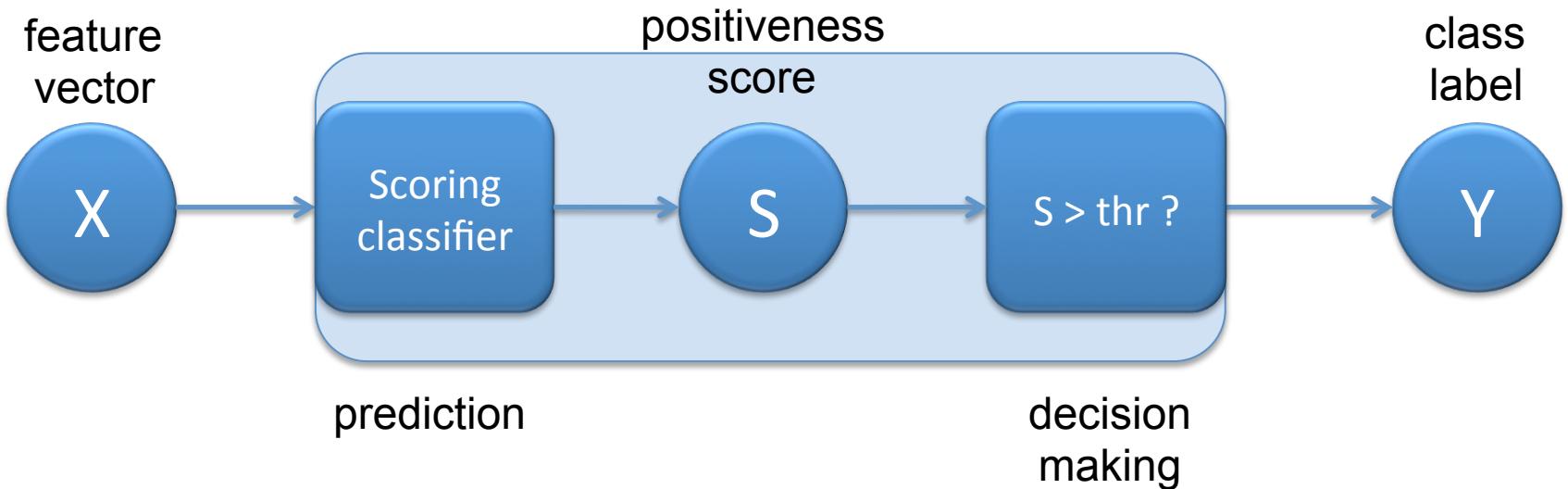


Prediction confidence of binary classifiers



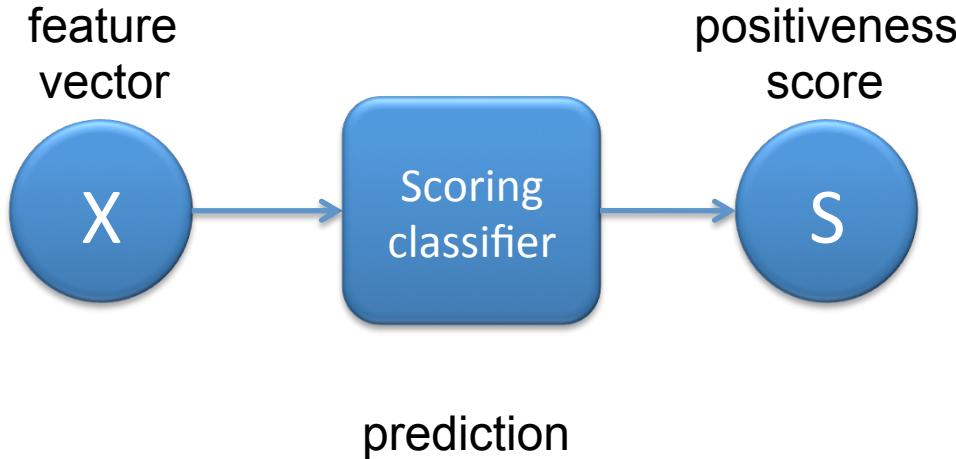
- Confidence use cases:
 - Abstaining classifier
 - Costs unknown during training
 - Human decision-maker
- Separation of prediction and decision-making

Score quantifies confidence in label



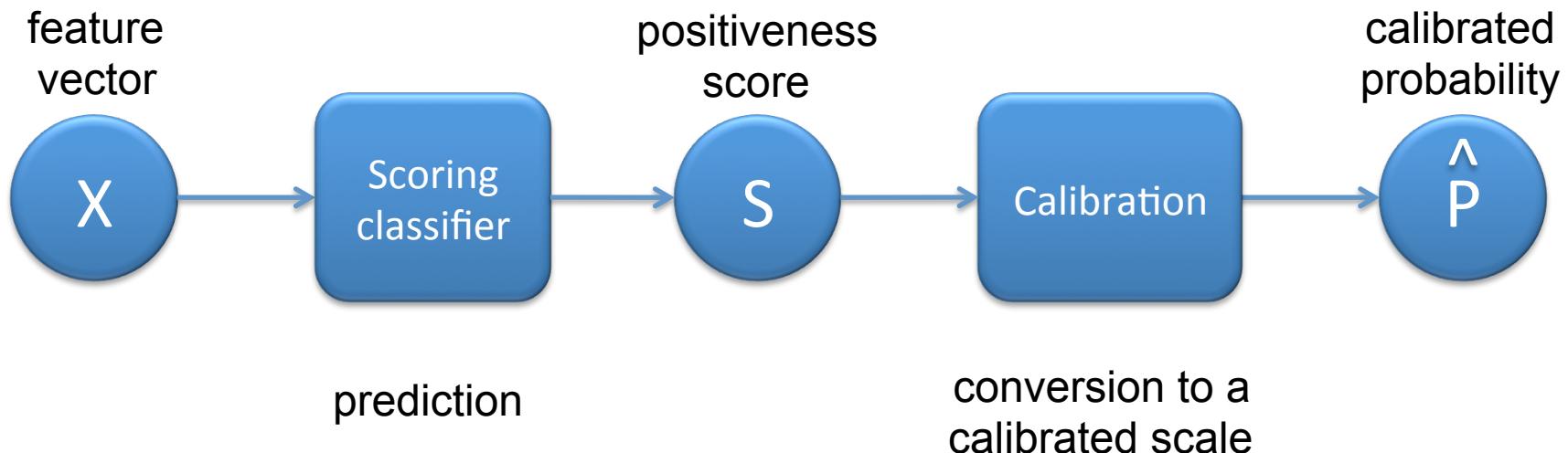
- Higher positiveness score S means more likely $Y=1$
 - SVM: signed distance to the decision boundary
 - Neural net: positive class probability from softmax
 - Naïve Bayes: positive class probability

Score quantifies confidence in label



- Higher positiveness score S means more likely $Y=1$
 - SVM: signed distance to the decision boundary
 - Neural net: positive class probability from softmax
 - Naïve Bayes: positive class probability
- Problem: the scale of S is model-dependent
 - How many positives among instances with $S=0.8$?

Calibration converts to a meaningful scale

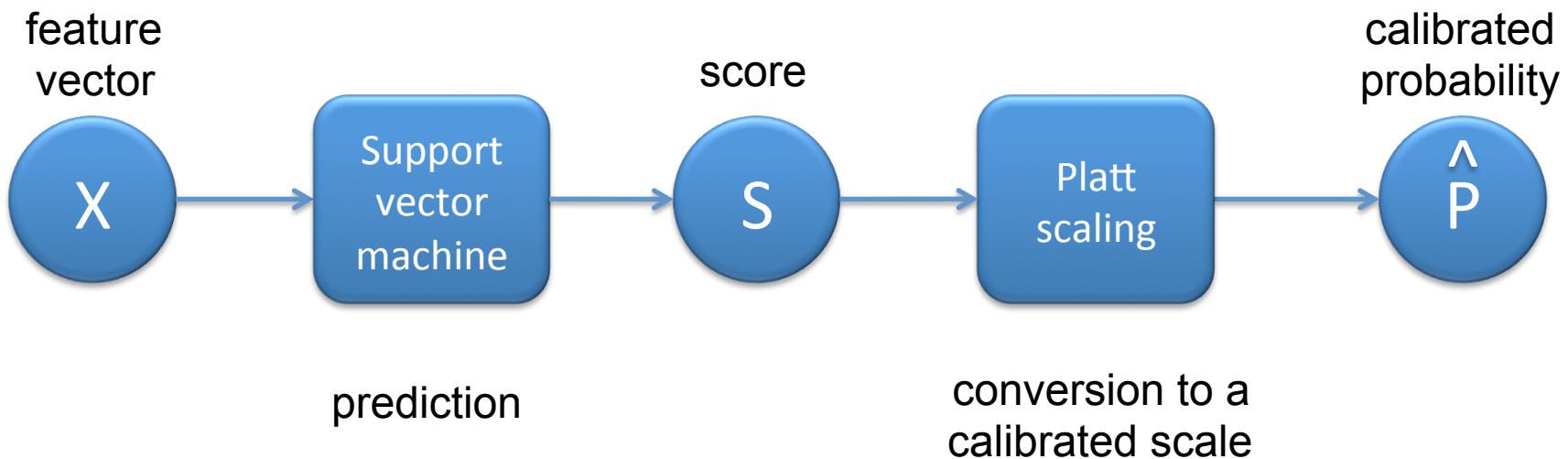


- Calibrated probabilities match with expected frequencies
 - E.g. among $\hat{P}=0.8$ there are 80% positive, 20% negative
 - Calibrated on the training data vs on the test data
- Perfect calibration map: $\hat{P} = P(Y=1 | S)$
- Bayes optimal: $\hat{P} = P(Y=1 | X)$

Logistic calibration (a.k.a. Platt scaling)

J. Platt (2000). Probabilities for SV machines.

In A. Smola, P. Bartlett, B. Schölkopf, D. Schuurmans (Eds.),
Advances in Large Margin Classifiers, pp.61–74, MIT Press.

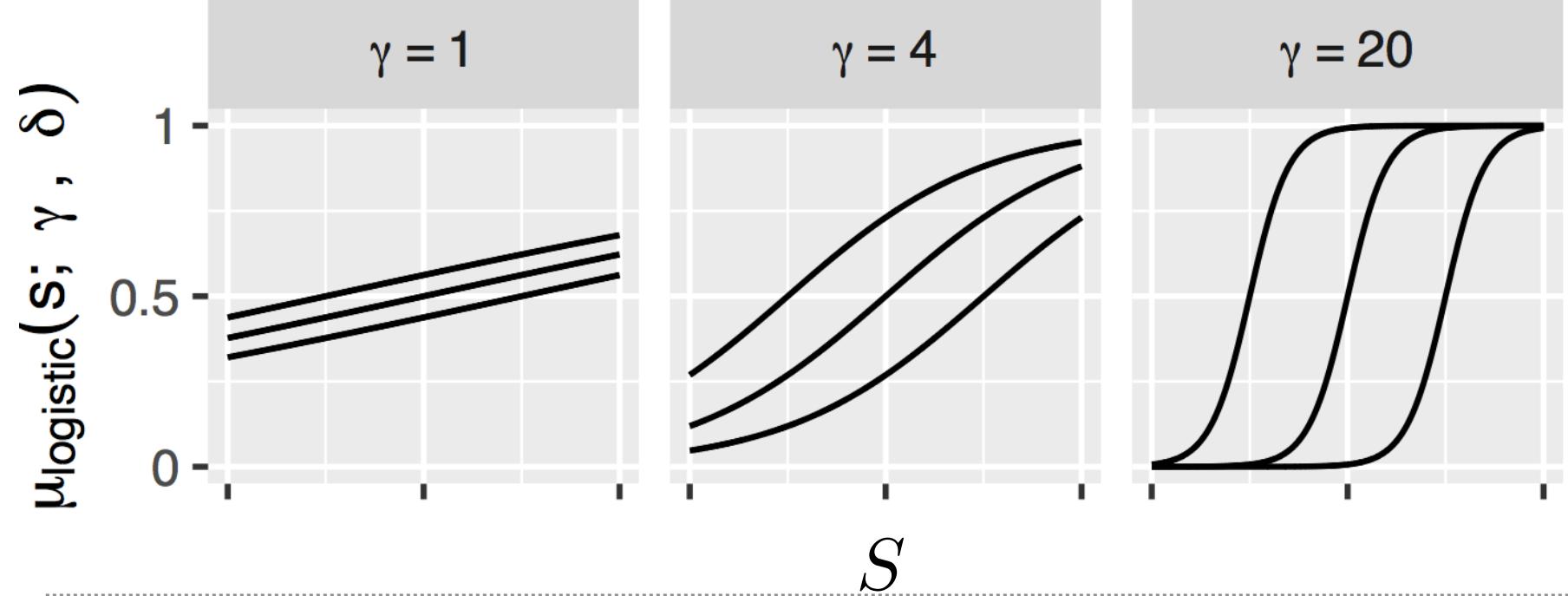


- Logistic sigmoid to convert $[-\infty, +\infty]$ to $[0, 1]$

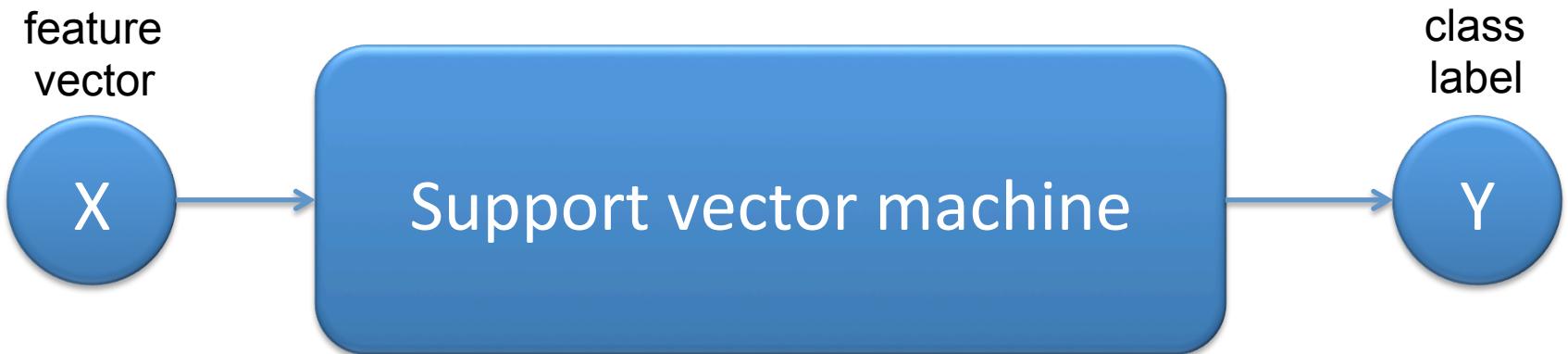
$$\mu_{logistic}(S; \gamma, \delta) = \frac{1}{1 + 1/(\exp(\gamma \cdot S + \delta))}$$

Logistic calibration map family

$$\mu_{logistic}(S; \gamma, \delta) = \frac{1}{1 + 1/(\exp(\gamma \cdot S + \delta))}$$



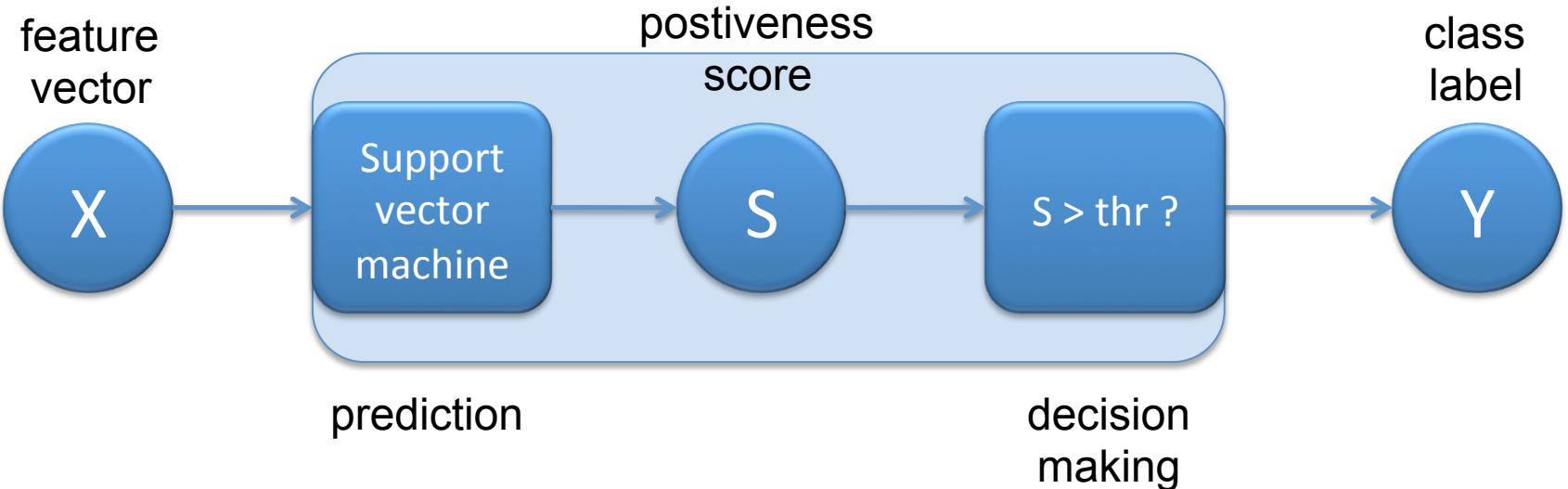
How to fit logistic calibration



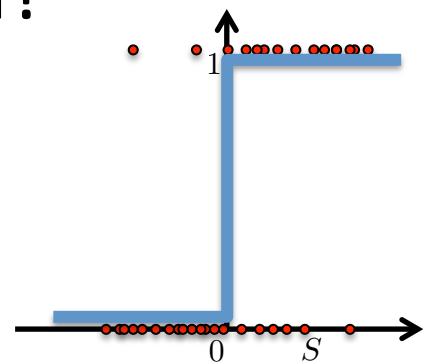
- How to get probabilities from SVM?



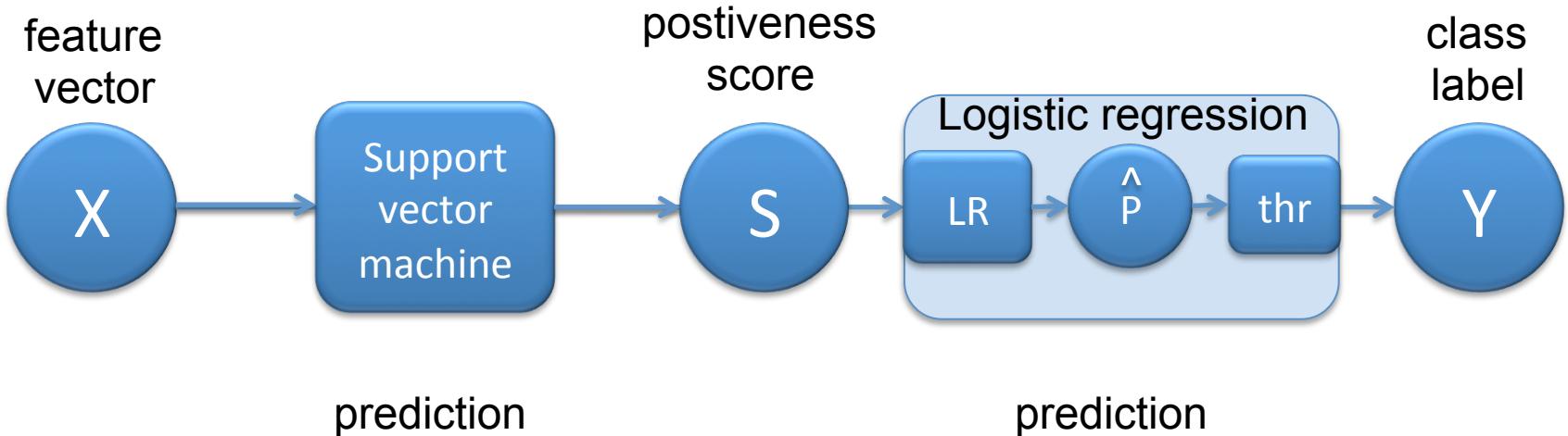
How to fit logistic calibration



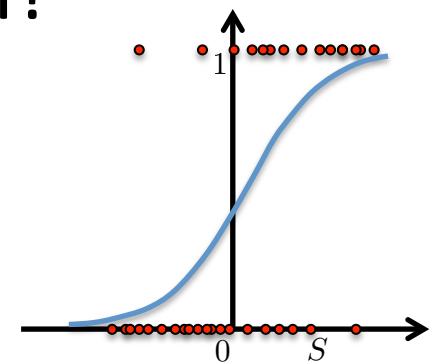
- How to get probabilities from SVM?
 1. Get the positiveness scores S



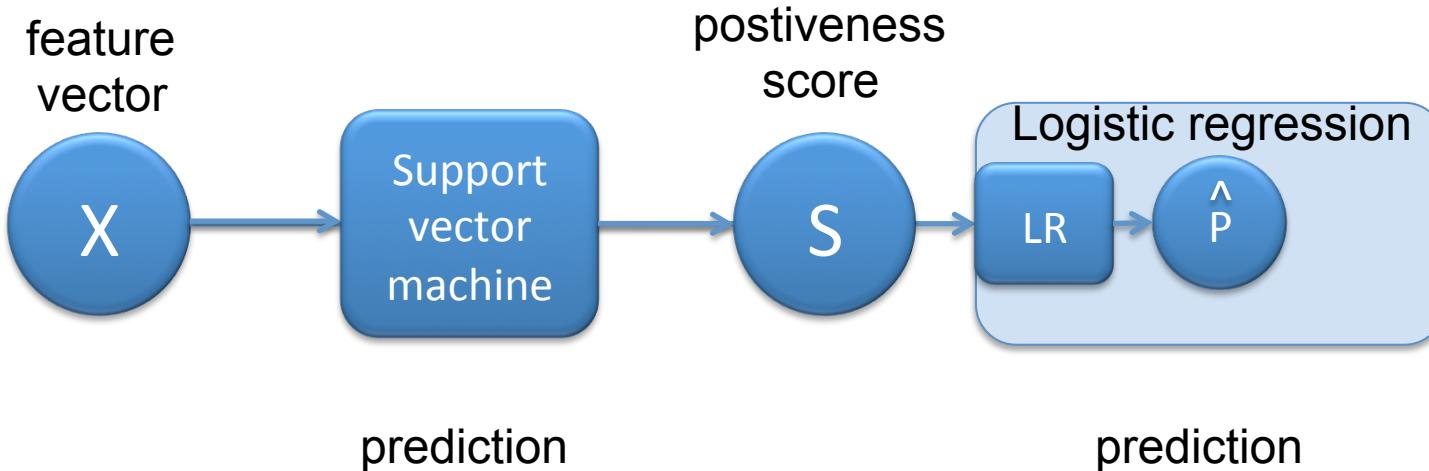
How to fit logistic calibration



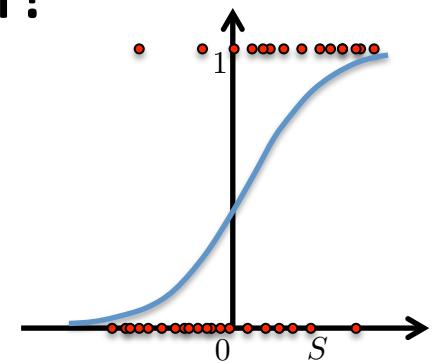
- How to get probabilities from SVM?
 1. Get the positiveness scores S
 2. Replace fixed step function by logistic regression fitting



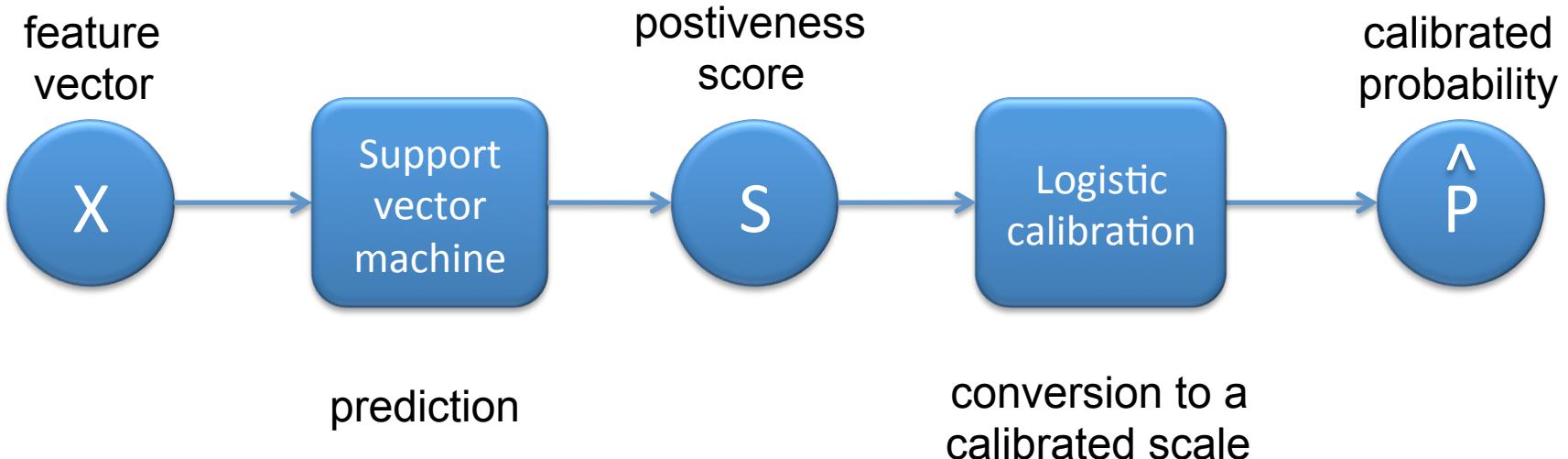
How to fit logistic calibration



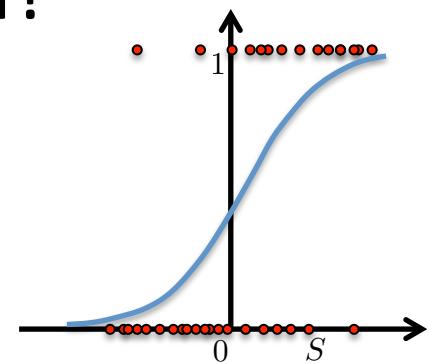
- How to get probabilities from SVM?
 1. Get the positiveness scores S
 2. Replace fixed step function by logistic regression fitting



How to fit logistic calibration



- How to get probabilities from SVM?
 1. Get the positiveness scores S
 2. Replace fixed step function by logistic regression fitting



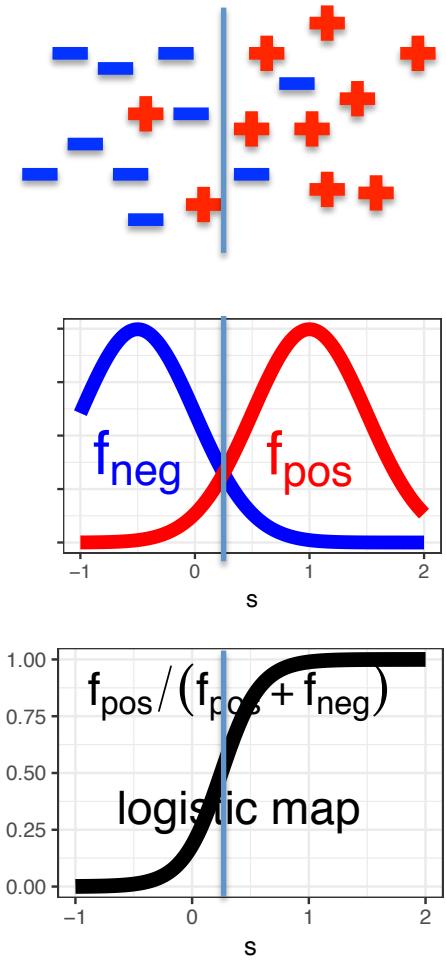
Logistic calibration is perfect if:

- Assume:

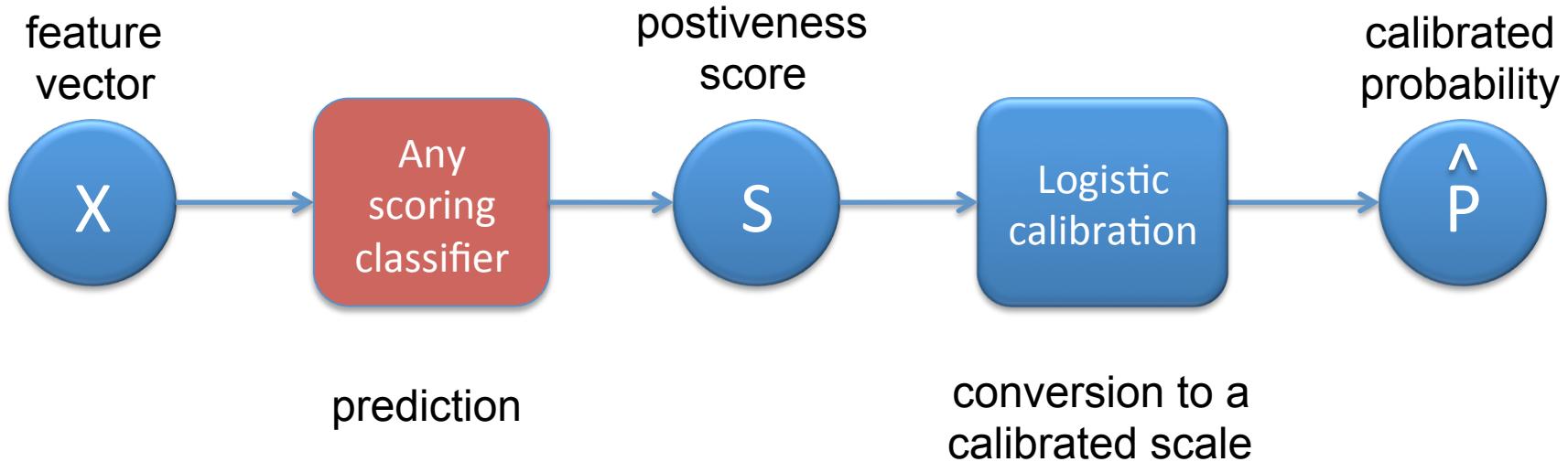
- $p(S | Y=+1) \sim \mathcal{N}(\mu_{pos}, \sigma^2)$
- $p(S | Y=-1) \sim \mathcal{N}(\mu_{neg}, \sigma^2)$
- $p(Y=+1) = p(Y=-1) = 0.5$

- Then:

$$\begin{aligned} p(Y=+1 | S) &= \frac{f_{pos}}{f_{pos} + f_{neg}} \\ &= \frac{1}{1 + 1/(\exp(\gamma \cdot S + \delta))} \end{aligned}$$

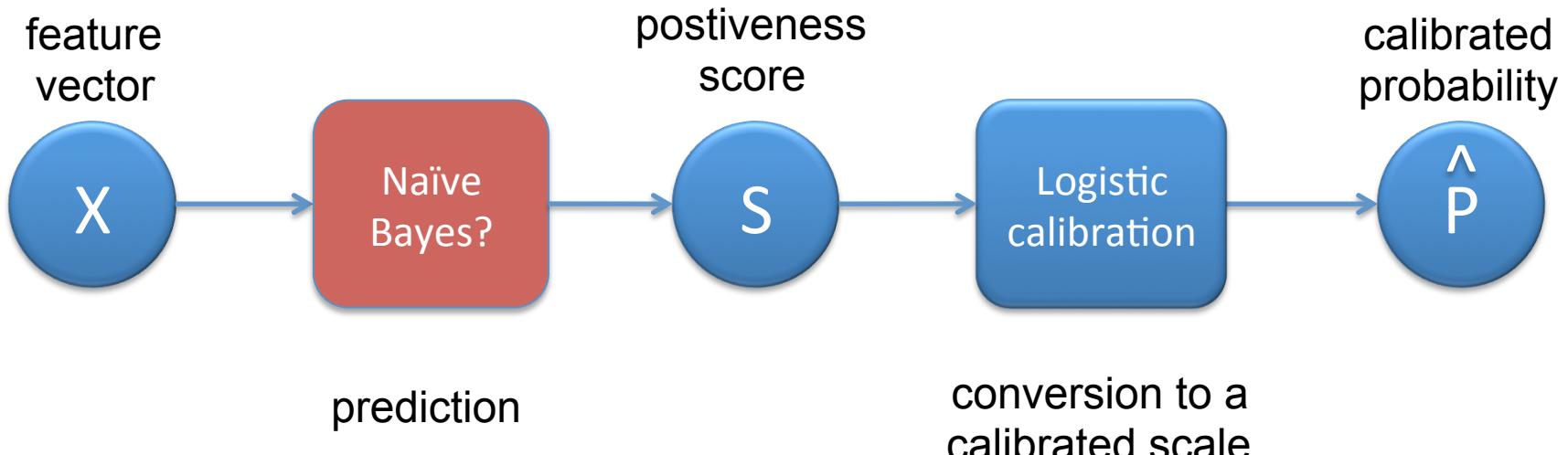


Logistic calibration is widespread

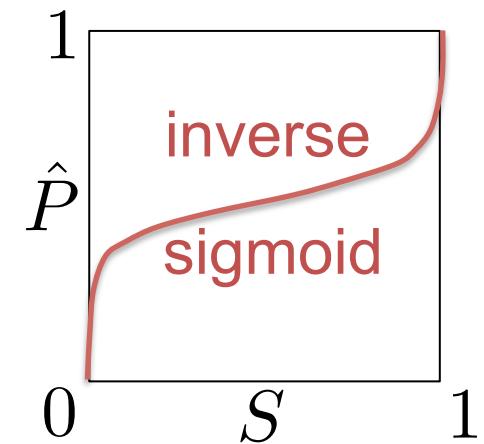


- If scores are already probabilities
 - Then Gaussian assumption questionable

Logistic calibration after Naïve Bayes?



- If scores are already probabilities
 - Then Gaussian assumption questionable
- Naïve Bayes pushes probabilities towards extremes
 - Particularly with many correlated features



Beta calibration from first principles:

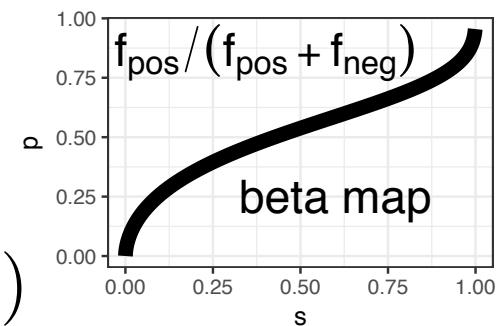
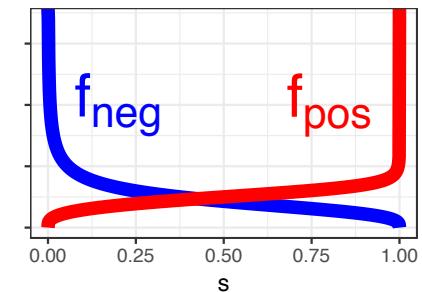
- Assume:

- $p(S | Y=+1) \sim Beta(\alpha_{pos}, \beta_{pos})$
- $p(S | Y=-1) \sim Beta(\alpha_{neg}, \beta_{neg})$
- $p(Y=+1) = p(Y=-1) = 0.5$

- Then:

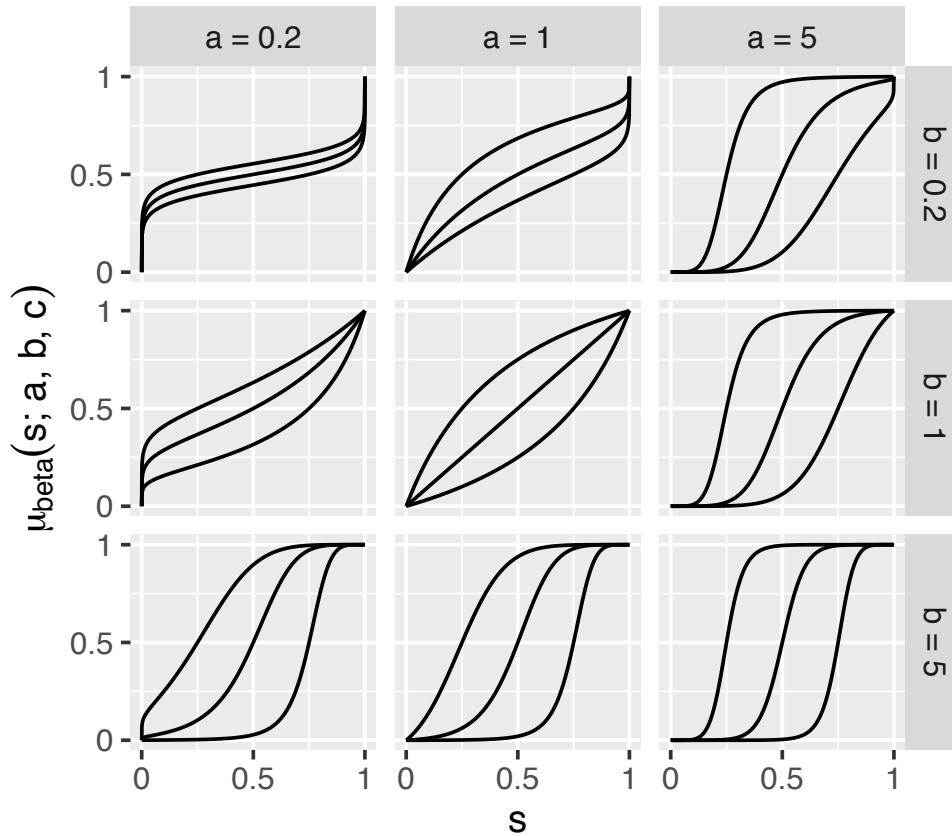
- $$p(Y=+1 | S) = \frac{f_{pos}}{f_{pos} + f_{neg}}$$

$$= \frac{1}{1 + 1/(e^c \frac{S^a}{(1-S)^b})} =: \mu_{beta}(S; a, b, c)$$



Beta calibration map family

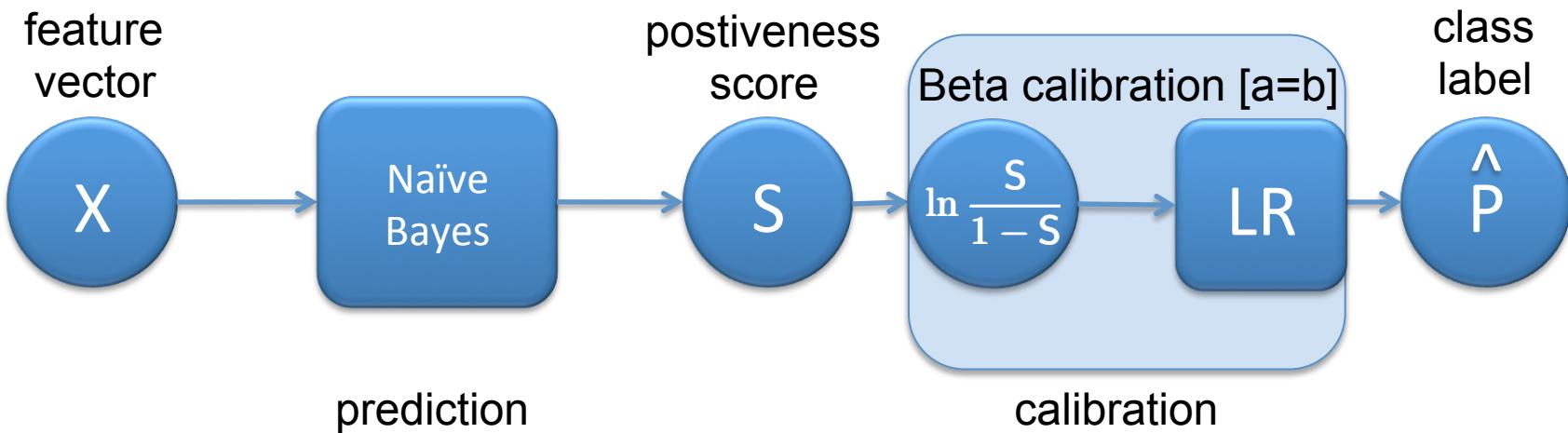
$$\mu_{beta}(S; a, b, c) = \frac{1}{1 + 1/(e^c \frac{S^a}{(1-S)^b})}$$



Logistic calibration vs Beta calibration

$$\mu_{logistic}(S; \gamma, \delta) = \frac{1}{1 + 1/(\exp(\gamma \cdot S + \delta))} \quad \mu_{beta}(S; a, b, c) = \frac{1}{1 + 1/(e^c \frac{S^a}{(1-S)^b})}$$

$$\mu_{logistic}\left(\ln \frac{S}{1 - S}; a, c\right) = \mu_{beta}(S; a, a, c)$$



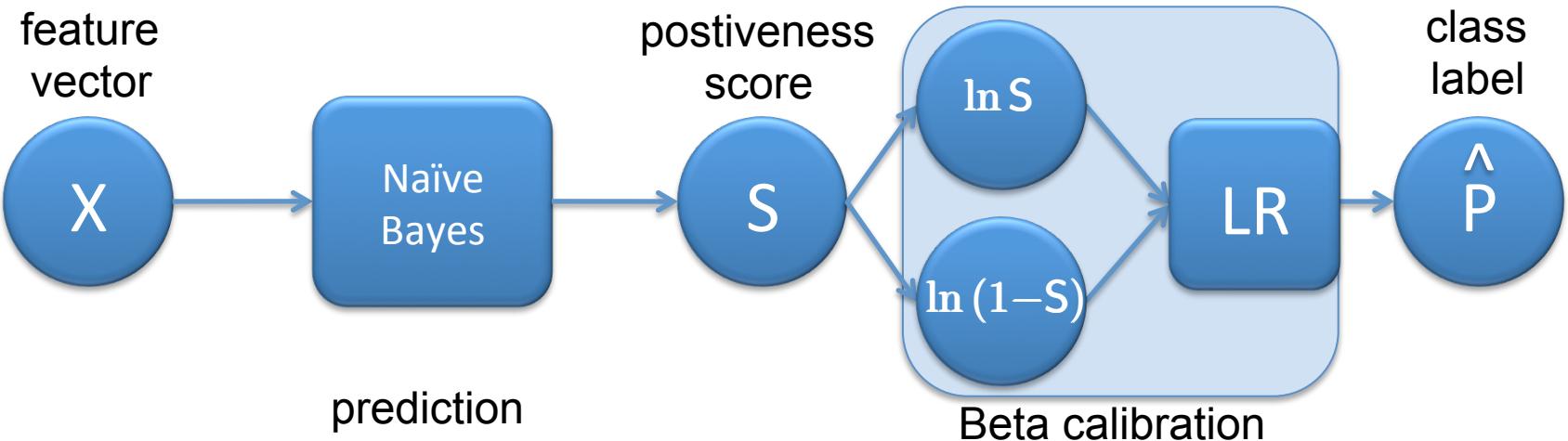
- 2-parameter beta calibration can be implemented as logistic calibration of the log odds

Beta calibration via logistic regression

$$\begin{aligned}\mu_{logreg}(S', S''; a, b, c) \\ = \frac{1}{1 + 1/(\exp(a \cdot S' + b \cdot S'' + c))}\end{aligned}$$

$$\mu_{beta}(S; a, b, c) = \frac{1}{1 + 1/(e^c \frac{S^a}{(1-S)^b})}$$

$$\mu_{logreg}(\ln S, -\ln(1 - S); a, b, c) = \mu_{beta}(S; a, b, c)$$



- 3-parameter beta calibration can be implemented as logistic regression on features $\ln S$ and $\ln(1-S)$

Experimental setup

- 41 UCI datasets
- 10x5-fold cross-validation
- Calibration with 3-fold internal cross-validation, predictions averaged [same as Platt 2000]
- Models:
 - NB : Naïve Bayes
 - Ada-O : Adaboost original
 - Ada-S : Adaboost in Scikit (SAMME)
- Evaluation measures: proper scoring rules
 - Brier score (squared error) $BS = \mathbb{E}[(\hat{P} - \mathbb{1}_{Y=1})^2]$
 - Log-loss (cross-entropy) $LL = \mathbb{E}[-\ln(\hat{P}) \mathbb{1}_{Y=1} - \ln(1 - \hat{P}) \mathbb{1}_{Y=-1}]$

Results (Adaboost-original, log-loss)

Table 2: Log-loss results with Ada-O classifier. Best results are marked in **bold** and subscripts indicate the ranks (before rounding to 3 decimal digits).

dataset	uncalibrated	beta	isotonic	logistic	landsat	0.041 ₂	0.040 ₁	0.043 ₃	0.053 ₄
abalone	0.612 ₂	0.612 ₁	0.626 ₄	0.614 ₃	letter	0.037 ₃	0.035 ₂	0.035 ₁	0.044 ₄
autos	0.427 ₄	0.283 ₁	0.416 ₃	0.287 ₂	libras-	0.589 ₄	0.100 ₁	0.184 ₃	0.112 ₂
balance	0.049 ₄	0.037 ₁	0.041 ₃	0.038 ₂	lung-ca	0.193 ₄	0.139 ₁	0.189 ₃	0.139 ₂
car	0.114 ₂	0.109 ₁	0.122 ₄	0.119 ₃	mfeat-k	0.062 ₄	0.027 ₁	0.048 ₃	0.032 ₂
clevela	0.496 ₃	0.417 ₁	0.523 ₄	0.429 ₂	mfeat-m	0.040 ₄	0.014 ₁	0.026 ₃	0.014 ₂
credit-	0.360 ₃	0.341 ₂	0.444 ₄	0.338 ₁	mfeat-z	0.095 ₄	0.032 ₁	0.063 ₃	0.039 ₂
dermato	0.036 ₄	0.019 ₁	0.035 ₃	0.021 ₂	mushroo	0.000 ₄	0.000 ₃	0.000 ₂	0.000 ₁
diabete	0.506 ₃	0.484 ₁	0.525 ₄	0.495 ₂	optdigi	0.035 ₂	0.033 ₁	0.044 ₄	0.040 ₃
ecoli	0.326 ₄	0.145 ₁	0.224 ₃	0.159 ₂	page-bl	0.093 ₂	0.089 ₁	0.098 ₃	0.109 ₄
flare	0.404 ₁	0.404 ₂	0.421 ₄	0.408 ₃	pendigi	0.017 ₂	0.017 ₁	0.023 ₄	0.021 ₃
german	0.511 ₃	0.506 ₂	0.538 ₄	0.506 ₁	scene-c	0.384 ₄	0.364 ₁	0.376 ₂	0.378 ₃
glass	0.696 ₄	0.489 ₂	0.560 ₃	0.485 ₁	segment	0.020 ₃	0.010 ₁	0.023 ₄	0.013 ₂
heart-s	0.570 ₄	0.427 ₁	0.557 ₃	0.443 ₂	shuttle	0.000 ₂	0.000 ₁	0.001 ₄	0.000 ₃
hepatit	0.805 ₄	0.392 ₁	0.431 ₃	0.411 ₂	sonar	0.941 ₄	0.404 ₁	0.486 ₃	0.440 ₂
horse	0.642 ₄	0.413 ₁	0.524 ₃	0.418 ₂	spambas	0.162 ₂	0.162 ₁	0.173 ₄	0.167 ₃
ionosph	0.381 ₄	0.203 ₁	0.296 ₃	0.227 ₂	tic-tac	0.379 ₄	0.333 ₁	0.340 ₃	0.339 ₂
iris	0.127 ₄	0.000 ₂	0.000 ₁	0.000 ₃	vehicle	0.077 ₂	0.068 ₁	0.131 ₄	0.077 ₃
					vowel	0.076 ₂	0.071 ₁	0.094 ₄	0.085 ₃
					wavefor	0.253 ₂	0.252 ₁	0.264 ₃	0.271 ₄
					wdbc	0.256 ₄	0.089 ₁	0.141 ₃	0.107 ₂
					wpbc	1.016 ₄	0.500 ₂	0.496 ₁	0.503 ₃
					yeast	0.510 ₂	0.509 ₁	0.543 ₄	0.514 ₃
					zoo	0.132 ₄	0.013 ₃	0.013 ₁	0.013 ₂
					rank	3.20	1.27	3.12	2.41

Key findings

- Beta calibration better than logistic calibration:

	Naïve Bayes	Adaboost-O	Adaboost-S
Brier score	++	++	+
Log-loss	++	++	+

- Beta calibration better than isotonic calibration in log-loss:

	Naïve Bayes	Adaboost-O	Adaboost-S
Brier score	+	+	-
Log-loss	++	++	++

- Full 3-parameter version better than 2-parameter versions (but often non-significantly)



Availability

- <https://betacal.github.io>
 - Python package
 - R package
- Or, implement your own as logistic regression on the dataset with two features: $\ln s$ and $\ln(1 - s)$



Next: Dirichlet calibration for multi-class

- Dirichlet distribution is the multi-class generalisation of beta distribution
- Beta calibration can also be generalised in a relatively straight-forward manner
- <https://dircal.github.io> (paper in preparation)
- Or, implement yourself as logistic regression on features $\ln p_1, \ln p_2, \dots, \ln p_k$ where k is the number of classes



Summary

- Scoring classifiers
 - Allow for separation of prediction and decision making
- Calibrated probabilities
 - Provide a standard scale of scores
- Logistic calibration
 - Is perfect when score conditional densities are Gaussian
- Beta calibration
 - Is perfect when score conditional densities are Beta
- Beta calibration improves logistic calibration if scores are in [0,1] range :
 - Well-founded, easily implemented, better performance
 - <https://betacal.github.io>

