Mathematical Tools for Neural and Cognitive Science 1

2

3

Fall semester, 2018

Probability & Statistics: Estimation, inference, model-fitting

Estimation of model parameters (outline)

- How do I compute an estimate? (mathematics vs. numerical optimization)
- How "good" are my estimates? (classical stats vs. simulation vs. resampling)
- How well does my model explain the data? Future data (prediction/generalization)? (classical stats vs. resampling)
- How do I compare two (or more) models? (classical stats vs. resampling)

The sample average $a(\vec{x}) = \frac{1}{N} \sum_{n=1}^{N} x_n$



- Most common common form of estimator
- Value of *a* converges to true mean *E*(*x*), for all reasonable distributions
- Variance of *a* converges to zero, as σ^2/N
- Distribution *p(a)* converges to a Gaussian (the "Central Limit Theorem")













Point Estimates

- Estimator: Any function of the data, intended to provide an estimate of the true value of a parameter
- Statistically-motivated estimators:
 - Maximum likelihood (ML):
 - Max a posteriori (MAP):
 - Bayes estimator:
 - Bayes least squares: (special case)

 $\hat{x}(\vec{d}) = \arg\min_{\hat{x}} \mathbf{E} \left(L(x - \hat{x}) | \vec{d} \right)$ $\hat{x}(\vec{d}) = \arg\min_{\hat{x}} \mathbf{E} \left((x - \hat{x})^2 | \vec{d} \right)$ $= \mathbf{E} \left(x | \vec{d} \right)$

9

Estimator quality: Bias & Variance

- Mean squared error = $bias^2 + variance$
- Bias is difficult to assess (requires knowing the "true" value). Variance is easier.
- Classical statistics generally aims for an unbiased estimator, with minimal variance ("MVUE").
- The MLE is *asymptotically* unbiased (under fairly general conditions), but this is only useful if
 - the likelihood model is correct
 - the optimum can be computed
 - you have lots of data
- More general view: estimation is about trading off bias and variance, through model selection, "regularization", or Bayesian priors...





























prior fair

0.4 0.6 x

1.5

) d 1

0.5

ob ob 0.2

(H)X)d

ob b



























Significance levels

- Are denoted by the Greek letter α .
- In principle, we can pick anything that we consider unlikely.
- In practice, the consensus is that a level of 0.05 or 1 in 20 is considered as unlikely enough to reject H₀ and accept the alternative.
- A level of 0.01 or 1 in 100 is considered "highly significant" or really unlikely.

Does NZT improve IQ scores or not?		28	
	Real	lity	
	Yes	No	
ïcant? Yes	Correct	Type I error α-error False alarm	
Signif No	Type II error β -error Miss	Correct	

28		

Test statistic

29

30

- We calculate how far the observed value of the sample average is away from its expected value.
- In units of standard error.
- In this case, the test statistic is

$$z = \frac{\overline{x} - \mu}{SE} = \frac{\overline{x} - \mu}{\sigma / \sqrt{N}}$$

• Compare to a distribution, in this case z or N(0,1)

Common misconceptions

Is "Statistically significant" a synonym for:

- Substantial
- Important
- Big
- Real

Does statistical significance gives the

- probability that the null hypothesis is true
- probability that the null hypothesis is false
- probability that the alternative hypothesis is true
- probability that the alternative hypothesis is false

Meaning of *p*-value. Meaning of CI.

Student's t-test

• σ not assumed known • Use $\sum_{n=1}^{N} (-1)^{2}$

$$s^{2} = \frac{\sum_{i=1}^{N} (x_{i} - \overline{x})^{2}}{N - 1}$$

• Why *N*-1? *s* is unbiased (unlike ML version), i.e., $E(s^2) = \sigma^2$

• Test statistic is $t = \frac{\overline{x} - \mu_0}{s / \sqrt{N}}$

- Compare to *t* distribution for CIs and NHST
- "Degrees of freedom" reduced by 1 to N-1







Many varieties of frequentist univariate tests

- χ^2 goodness of fit
- χ^2 test of independence
- test a variance using χ^2
- *F* to compare variances (as a ratio)
- Nonparametric tests (e.g., sign, rank-order, etc.)

Bootstrapping

35

- "The Baron had fallen to the bottom of a deep lake. Just when it looked like all was lost, he thought to pick himself up by his own bootstraps" [Adventures of Baron von Munchausen, by Rudolph Erich Raspe]
- A (**re**)**sampling** method for computing estimator distribution (incl. stdev error bars or confidence intervals)
- Idea: instead of running experiment multiple times, resample (with replacement) from the *existing* data. Compute an estimate from each of these "bootstrapped" data sets.



For strokes, the ratio of rates is

$$\widehat{\theta} = \frac{119/11037}{98/11034} = 1.21. \tag{1.4}$$

It now looks like taking aspirin is actually harmful. However the interval for the true stroke ratio θ turns out to be

$$93 < \theta < 1.59$$
 (1.5)

with 95% confidence. This includes the neutral value $\theta = 1$, at which aspirin would be no better or worse than placebo vis-à-vis strokes. In the language of statistical hypothesis testing, aspirin was found to be significantly beneficial for preventing heart attacks, but not significantly harmful for causing strokes.

[Efron & Tibshirani '98]





- Bayes estimator:

























































Decision/classification in multiple dimensions

- Data-driven:
 - Fisher Linear Discriminant (FLD) maximize d'
 - Support Vector Machine (SVM) maximize margin
- Statistical:
 - ML/MAP/Bayes under a probabilistic model
 - e.g.: Gaussian, equal covariance (same as FLD)
 - e.g.: Gaussian, unequal covariance (QDA)
- Examples:
 - Visual gender identification
 - Neural population decoding



































 (θ_2)







72



Take N independent samples from the distribution, these act like draws from N independent, identically distributed (IID) RVs:

 $X_1, X_2, \cdots X_N$

Continuous/Gaussian: Localization

The *N* independent samples are $x_1, x_2, \dots x_N$

ML estimates are

 $\hat{\sigma}^2 = \frac{\sum_{i=1}^{N} (x_i - \overline{x})^2}{N - 1}$

 $\hat{\mu} = \frac{\sum_{i=1}^{N} x_i}{N}$

Continuous/Gaussian: Localization

MAP estimates of the mean are based on the posterior, a product of Gaussians (assuming a Gaussian) prior. Thus there is shrinkage toward the prior.

Model comparison for hypotheses about the mean (variance assumed known) are similar to the binomial example.







Product of Gaussian distributions is Gaussian







































































Taxonomy of model-fitting errors

- Optimization failures (e.g., local minima) [convex relaxation, test with simulations]
- Overfitting (too many params, not enough data) [use cross-validation to select complexity, or to control regularization]
- Experimental variability (due to finite/noisy measurements) [use math/distributional assumptions, or simulations, or bootstrapping]
- Model failures



98	





Cross-validation

100

A resampling method for constraining a model. Widely used to identify/avoid over-fitting.

 (1) Randomly partition data into a "training" set, and a "test" set.
(2) Fit model to training set. Measure error on test set.
(3) Repeat (many times)

(4) Choose model that minimizes the cross-validated (test) error



Using cross-validation to select the degree of a polynomial model:







L1 regularization

(a.k.a. least absolute shrinkage and selection operator - LASSO)















K-Means

108

• Estimating cluster assignments: given class centers, assign each point to closest one.





• Estimating cluster parameters: given assignments, reestimate the centroid of each cluster.









ML for discrete mixture of Gaussians

$$p(\vec{x}_n | a_{nk}, \vec{\mu}_k, \Lambda_k) \propto \sum_k \frac{a_{nk}}{\sqrt{|\Lambda_k|}} e^{-(\vec{x}_n - \vec{\mu}_k)^T \Lambda_k^{-1} (\vec{x}_n - \vec{\mu}_k)/2}$$
$$a_{nk} = \text{assignment } probability$$

 $\{\vec{\mu}_k, \Lambda_k\} = \text{mean/covariance of class } n$

Intuition: alternate between maximizing these two sets of variables ("coordinate descent")













