

Foundations of

Machine Learning















Distributions & Data · Instance linput space X · E.g. Loan apps, images, med. records,... · Outcome/output space Y · E.g. loan status, cats, d'agnosis,... · Examples <x,y>:xeX,geY · Distribution/population P over possible <x,y7 • All we see is a sample S: S= E < x1, y17, < x2, g27, ..., < x1, yn73 drawn from P







Error: Train & True · Training eror of hon S: $\hat{\varepsilon}_{s}(h) = \hat{\varepsilon}(h) \stackrel{\text{def}}{=} \stackrel{\text{def}}{=} \stackrel{\text{def}}{=} I[h(x_{i}) \neq y_{i}]$ ("indicator function") = fraction of mistakes of hon S LcLassification setting) • True/fest error of h on P: Ep(h)=E(h)=E[I]h(x)=y] (x,y>-P = probability h makes a mistake on xx,y7~P





"simple" models

· Choose hEH where H"simple"









true error E(h)

Standard ML Workflow · Gather sample Show P · Choose/design model class H · Use algo/heuristic to find ht H with small \$(h) · Estimate Elh) on new dota also from P · (Repeat...) What justifics this methodology?





P is some pop./distribution
over 4x, y7 pairs
Uchave sample

























· If we use squared error:











OK. But again,

why is this a good idea?



ML Rescarch

· Design of notural lexpressive H

· Design of fost algos for

(approx) minimizing Elh) in H









Fairness, privacy,

explainability, safety,

robustness...

