

Intuition

How about if we inferred **properties** of a subset of the training inputs...

...but not the whole class?

In a nutshell: given a **gender** classifier, infer **race** of
people in Bob's photos





In a nutshell: given a gender classifier, inference

training inputs. . .

...but not the whole class?

How about if we inferred **properties** of a subset of the

people in Bob's photos

